

## AN EFFICIENT MULTICORE IMPLEMENTATION OF A NOVEL HSS-STRUCTURED MULTIFRONTAL SOLVER USING RANDOMIZED SAMPLING\*

PIETER GHYSELS<sup>†</sup>, XIAOYE S. LI<sup>†</sup>, FRANÇOIS-HENRY ROUET<sup>†</sup>,  
SAMUEL WILLIAMS<sup>‡</sup>, AND ARTEM NAPOV<sup>‡</sup>

**Abstract.** We present a sparse linear system solver that is based on a multifrontal variant of Gaussian elimination and exploits low-rank approximation of the resulting dense frontal matrices. We use hierarchically semiseparable (HSS) matrices, which have low-rank off-diagonal blocks, to approximate the frontal matrices. For HSS matrix construction, a randomized sampling algorithm is used together with interpolative decompositions. The combination of the randomized compression with a fast ULV HSS factorization leads to a solver with lower computational complexity than the standard multifrontal method for many applications, resulting in speedups up to sevenfold for problems in our test suite. The implementation targets many-core systems by using task parallelism with dynamic runtime scheduling. Numerical experiments show performance improvements over state-of-the-art sparse direct solvers. The implementation achieves high performance and good scalability on a range of modern shared memory parallel systems, including the Intel Xeon Phi (MIC). The code is part of a software package called STRUMPACK (STRUctured Matrices PACKage), which also has a distributed memory component for dense rank-structured matrices.

**Key words.** sparse Gaussian elimination, multifrontal method, HSS matrices, parallel algorithm

**AMS subject classifications.** 65F05, 65F50, 97N80

**DOI.** 10.1137/15M1010117

**1. Introduction.** Solving large linear systems efficiently on modern hardware is an important requirement for many engineering high performance computing codes. For a wide range of applications, like those using finite element, finite difference, or finite volume discretizations of partial differential equations (PDEs), the resulting linear system is extremely sparse. Fast solution methods exploit this sparsity, but also arrange the computations in such a way that most of the computational work is done on smaller dense submatrices. The reason for this is that operations on dense matrices can be implemented very efficiently on modern hardware. The multifrontal method [23, 38] is an example of a sparse direct solver where most of the work is done on dense, so-called frontal, matrices. Unfortunately, dense linear algebra operations, for instance LU decomposition, require  $\mathcal{O}(N^3)$  operations, where  $N$  is the matrix

---

\*Received by the editors February 25, 2015; accepted for publication (in revised form) May 12, 2016; published electronically October 27, 2016. Partial support for this work was provided through Scientific Discovery through Advanced Computing (SciDAC) program funded by U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research (and Basic Energy Sciences/Biological and Environmental Research/High Energy Physics/Fusion Energy Sciences/Nuclear Physics). This research used resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy under contract DE-AC02-05CH11231. This manuscript has been authored by an author at Lawrence Berkeley National Laboratory under contract DE-AC02-05CH11231 with the U.S. Department of Energy. The U.S. Government retains, and the publisher, by accepting this article for publication, acknowledges, that the U.S. Government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for U.S. Government purposes.

<http://www.siam.org/journals/sisc/38-5/M101011.html>

<sup>†</sup>Lawrence Berkeley National Laboratory, Berkeley, CA 94720 (pghysels@lbl.gov, xsli@lbl.gov, frouet@lbl.gov, swilliams@lbl.gov).

<sup>‡</sup>Université Libre de Bruxelles, B-1050 Brussels, 1150, Belgium (anapov@ulb.ac.be).

dimension. In a multifrontal solver these dense operations end up being the bottleneck. However, it has been observed that for many applications the dense frontal matrices have some kind of low-rank structure [17].

In [55], a rank-structured multifrontal method is presented in which the larger frontal matrices are approximated by hierarchically semiseparable (HSS) [50] matrices. For certain model problems, this leads to a solver, or preconditioner, with linear or almost linear complexity in the total number of degrees of freedom in the sparse linear system. Here, we present an efficient implementation of a slightly modified version of the algorithm presented in [55]. The algorithm in [55] handles only symmetric positive definite systems, while the code presented here targets general nonsymmetric nonsingular matrices. For HSS compression, a randomized sampling algorithm from [39] is used. Earlier HSS construction methods (see [56]) cost at least  $\mathcal{O}(N^2)$ , whereas the randomized method in combination with a fast matrix-vector product has a linear or almost linear complexity, depending on the rank-structure of the frontal matrix.

An important concept used in the randomized compression algorithm is the interpolative or skeleton decomposition [19]. Use of this decomposition leads to a special structure of the HSS generator matrices (see (2.12)). The HSS factorization used in [55] for symmetric matrices, and in [57] for nonsymmetric matrices, exploits this special structure in a so-called ULV-like decomposition. In this paper, the ULV decomposition from [57] is used.

The HSS format is a subclass of a more general type of hierarchical rank-structured matrices called  $\mathcal{H}$ -matrices [13]. HSS matrices are similar to  $\mathcal{H}^2$ -matrices, another subclass of  $\mathcal{H}$ -matrices, in the sense that both formats have the special property that the generators are hierarchically nested (see (2.3) for what this means for HSS). This is typically not the case in the more general  $\mathcal{H}$ , the sequentially semiseparable (SSS) [50], or the hierarchically off-diagonal low-rank (HODLR) [3] formats (all of which are  $\mathcal{H}$ -matrices). In HSS and HODLR only off-diagonal blocks are approximated as low-rank, whereas  $\mathcal{H}$  and  $\mathcal{H}^2$  allow more freedom in the partitioning. In [4], a flat tiled block low-rank (BLR) format is used to approximate dense frontal matrices in the multifrontal solver MUMPS [6], while in other recent work [8] HODLR has also been proposed to accelerate a multifrontal solver. Both HSS and HODLR use similar hierarchical off-diagonal partitioning, but HSS further exploits the hierarchically nested bases' structure, which can lead to an asymptotically faster factorization algorithm for some matrices.

Furthermore, thanks to the randomized HSS construction, our solver is also fully structured (compared to partially structured approaches where only part of the frontal matrix is compressed; see [51]), and the larger frontal matrices are never explicitly formed as dense matrices.

Achieving high performance on multi-/many-core architectures can be challenging, but it has been demonstrated by many authors now that dynamic scheduling of fine-grained tasks represented by a directed acyclic graph (DAG) can lead to good performance for a range of codes. This approach was used successfully in the dense linear algebra libraries PLASMA and MAGMA [2], and more recently it has become clear that it is also a convenient and efficient strategy for sparse direct solvers. For instance, in [35] the PaStiX solver [30] is modified to use two different generic DAG schedulers (PaRSEC [14] and StarPU [9]). In [1] StarPU is used in a multifrontal QR solver. In [33], OpenMP tasks are submitted recursively for work on different frontal matrices, while parallelism inside the frontal matrices is also exploited with OpenMP tasks but with a manual resolution of intertask dependencies. The sparse Cholesky

solver HSL\_MA87 [31] uses a custom DAG scheduler implemented in OpenMP. Just as sparse direct solvers, hierarchical matrix algorithms also benefit from task parallelism: Kriemann [34] uses a DAG to schedule fine-grained tasks to perform  $\mathcal{H}$ -matrix LU factorization. Our implementation uses OpenMP task scheduling, but since most of the tasks are generated recursively, a DAG is never explicitly constructed.

The main contribution of this work is the development of a robust and efficient code for the solution of general sparse linear systems, with a specific emphasis on systems from the discretization of PDEs. Our work addresses various implementation issues, the most important being the use of an adaptive HSS construction scheme (section 2.2.1), based on the randomized sampling method [39]. Rather than assuming that the maximum rank in the HSS submatrices is known a priori, it is computed in an adaptive way during the HSS compression. Other implementation techniques such as fast extraction of elements from an HSS structure (section 4.5) are also indispensable to make the code robust and usable as a black-box solver. The code achieves high performance and good scalability on a range of modern multi-/many-core architectures like Intel Xeon and Intel Xeon Phi (MIC), due to runtime scheduled task parallelism, using OpenMP.<sup>1</sup> The exclusive use of task parallelism avoids expensive interthread synchronization and leads to a very scalable code. This is the first parallel algebraic sparse solver with fully structured HSS low-rank approximation. The code is made publicly available with a BSD license as part of a package called STRUMPACK<sup>2</sup> (STRUctured Matrices PACKage). STRUMPACK also has a dense distributed memory component; see [47].

This work advances the field significantly on several fronts.

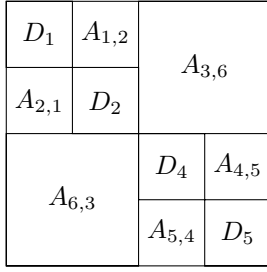
- Wang et al. [52] presented the first parallel multifrontal code with HSS embedding, called Hsolver. However, two shortcomings prevent it from being widely adopted: it is restricted to matrices arising from the discretization of regular meshes, and it is only partially structured due to the hurdle of the extend-add of HSS update matrices (see section 4). Our new code, on the other hand, is a purely algebraic solver for general sparse linear systems and is fully structured, mitigating the HSS extend-add obstacle thanks to the randomized sampling technique (see section 4.3).
- Napov and Li developed a purely algebraic sparse solver with HSS embedding [42], but it is only sequential and their experiments did not include the randomized sampling HSS compression. Our new code is parallel and we show detailed results with randomized sampling.

In future work, building on the current paper and on the distributed HSS code developed in [47], we intend to develop a distributed memory algebraic sparse solver with HSS compression.

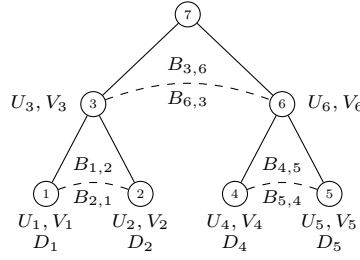
The rest of this paper is outlined as follows. Some required background on HSS is briefly presented in section 2. First, in section 2.1, the HSS rank-structured format is described. Next, the fast randomized sampling HSS construction [39] and the ULV decomposition [57] are discussed in sections 2.2 and 2.3, respectively. Section 3 describes multifrontal sparse LU decomposition. In section 4 we discuss how HSS matrices can be incorporated into a multifrontal solver. Section 5 explains various aspects of the actual implementation. In section 6 we present experimental results that illustrate numerical and performance aspects of the code. Finally, section 7 has some concluding remarks and an outlook on planned future work.

<sup>1</sup><http://openmp.org/wp/openmp-specifications/>

<sup>2</sup><http://portal.nersc.gov/project/sparse/strumpack/>



(a) HSS partitioning of a square matrix.



(b) Tree corresponding to the HSS partition.

For a leaf node,  $U_\tau = \hat{U}_\tau$  and  $V_\tau = \hat{V}_\tau$ .

FIG. 1. Illustration of an HSS partitioning of a square matrix. Diagonal blocks are partitioned recursively. Figure (b) shows the tree, using postordering, corresponding to the partitioning in (a), and it illustrates the basis matrices stored in the nodes of the tree.

**2. HSS: Hierarchically semiseparable matrices.** This section briefly introduces hierarchically semiseparable (HSS) matrices, mostly following the notation from [39]. HSS is a data-sparse matrix representation which is part of the more general class of  $\mathcal{H}$ -matrices and, more specifically,  $\mathcal{H}^2$ -matrices.

**2.1. Overview of the HSS matrix format.** The following notation is used: “,” is MATLAB-like notation for all indices in the range, “\*” denotes complex conjugation, “ $\#I_\tau$ ” is the number of elements in index set  $I_\tau = \{i_1, i_2, \dots, i_n\}$ , and  $R_\tau = R(I_\tau, :)$  is the matrix consisting of only the rows  $I_\tau$  of matrix  $R$ .

Consider a square matrix  $A \in \mathbb{C}^{N \times N}$  with an index set  $I_A = \{1, \dots, N\}$  associated with it. Let  $\mathcal{T}$  be a postordered binary tree, meaning that children in the tree are numbered before their parent. Each node  $\tau$  of the tree is associated with a contiguous subset  $t_\tau \subset I_A$ . For two siblings in the tree,  $\nu_1$  and  $\nu_2$ , children of  $\tau$ , it holds that  $t_{\nu_1} \cup t_{\nu_2} = t_\tau$  and  $t_{\nu_1} \cap t_{\nu_2} = \emptyset$ . Furthermore,  $\cup_{\tau=\text{leaf}(\mathcal{T})} t_\tau = t_{\text{root}(\mathcal{T})} = I_A$ . The same tree  $\mathcal{T}$  is used for the rows and the columns of  $A$ , and only diagonal blocks are partitioned. An example of the resulting matrix partitioning is given in Figure 1a, and the corresponding tree is shown in Figure 1b.

The diagonal blocks of  $A$ , denoted  $D_\tau$ , are stored as dense matrices in the leaves  $\tau$  of the tree  $\mathcal{T}$

$$(2.1) \quad D_\tau = A(I_\tau, I_\tau).$$

The off-diagonal blocks  $A_{\nu_1, \nu_2} = A(I_{\nu_1}, I_{\nu_2})$ , where  $\nu_1$  and  $\nu_2$  denote two siblings in the tree, are factored (approximately) as

$$(2.2) \quad A_{\nu_1, \nu_2} \approx \hat{U}_{\nu_1} B_{\nu_1, \nu_2} (\hat{V}_{\nu_2})^*.$$

The matrices  $\hat{U}_{\nu_1}$  and  $\hat{V}_{\nu_2}$ , which form bases for the column and row spaces of  $A_{\nu_1, \nu_2}$ , are typically tall and skinny, with  $\hat{U}_{\nu_1}$  having  $\#I_{\nu_1}$  rows and  $r_{\nu_1}^r$  (column-rank) columns,  $\hat{V}_{\nu_2}$  has  $\#I_{\nu_2}$  rows and  $r_{\nu_2}^c$  (row-rank) columns and hence  $B_{\nu_1, \nu_2}$  is  $r_{\nu_1}^r \times r_{\nu_2}^c$ . The HSS-rank  $r$  of matrix  $A$  is defined as the maximum of  $r_\tau^r$  and  $r_\tau^c$  over all off-diagonal blocks, where typically  $r \ll N$ . The matrices  $B_{\nu_1, \nu_2}$  and  $B_{\nu_2, \nu_1}$  are stored in the parent of  $\nu_1$  and  $\nu_2$ . For a nonleaf node  $\tau$  with children  $\nu_1$  and  $\nu_2$ ,

the basis matrices  $\hat{U}_\tau$  and  $\hat{V}_\tau$  are not stored directly since they can be represented hierarchically as

$$(2.3) \quad \hat{U}_\tau = \begin{bmatrix} \hat{U}_{\nu_1} & 0 \\ 0 & \hat{U}_{\nu_2} \end{bmatrix} U_\tau \quad \text{and} \quad \hat{V}_\tau = \begin{bmatrix} \hat{V}_{\nu_1} & 0 \\ 0 & \hat{V}_{\nu_2} \end{bmatrix} V_\tau.$$

Note that for a leaf node  $\hat{U}_\tau = U_\tau$  and  $\hat{V}_\tau = V_\tau$ . Hence, every node  $\tau$  with children  $\nu_1$  and  $\nu_2$ , except for the root node, keeps matrices  $U_\tau$  and  $V_\tau$  instead of the typically larger  $\hat{U}_\tau$  and  $\hat{V}_\tau$ . The example from Figure 1a can be written out explicitly as

$$(2.4) \quad A = \begin{bmatrix} D_1 & U_1 B_{1,2} V_2^* & \begin{bmatrix} U_1 & 0 \\ 0 & U_2 \end{bmatrix} U_3 B_{3,6} V_6^* \begin{bmatrix} V_4^* & 0 \\ 0 & V_5^* \end{bmatrix} \\ U_2 B_{2,1} V_1^* & D_2 & \begin{bmatrix} U_1 & 0 \\ 0 & U_2 \end{bmatrix} U_3 B_{3,6} V_6^* \begin{bmatrix} V_4^* & 0 \\ 0 & V_5^* \end{bmatrix} \\ \begin{bmatrix} U_4 & 0 \\ 0 & U_5 \end{bmatrix} U_6 B_{6,3} V_3^* \begin{bmatrix} V_1^* & 0 \\ 0 & V_2^* \end{bmatrix} & D_4 & U_4 B_{4,5} V_5^* \\ \begin{bmatrix} U_4 & 0 \\ 0 & U_5 \end{bmatrix} U_6 B_{6,3} V_3^* \begin{bmatrix} V_1^* & 0 \\ 0 & V_2^* \end{bmatrix} & U_5 B_{5,4} V_4^* & D_5 \end{bmatrix}.$$

The storage requirements for an HSS matrix are  $\mathcal{O}(rN)$ . Construction of the HSS generators will be discussed in the next section. Once an HSS representation of a matrix is available, it can be used to perform matrix-vector multiplication in  $\mathcal{O}(rN)$  operations compared to  $\mathcal{O}(N^2)$  for classical dense matrix-vector multiplication; see [39, 47].

**2.2. Fast HSS construction through randomized sampling.** In [39] Martinsson presents a randomized sampling algorithm for the efficient construction of an HSS representation of a matrix  $A$ . Note that the same technique was also used by Xia et al. in [57, 55] for HSS compression in a multifrontal solver. The main advantage of this approach is that it does not need explicit access to all elements of  $A$ , but only needs a fast matrix-vector routine and selected elements from  $A$ . The matrix  $A$  never needs to be formed explicitly as a dense matrix, and this allows us to save memory. The overall complexity of the algorithm is  $\mathcal{O}(Nr^2)$ , with  $r$  the HSS-rank of  $A$ , provided that a fast ( $\mathcal{O}(N)$ ) matrix-vector product is available. This section briefly presents the randomized compression algorithm. For a more in-depth discussion, see [47, 39].

Suppose the HSS-rank  $r$  is known a priori and  $R \in \mathbb{C}^{N \times d}$  is a tall and skinny random matrix with  $d = r + p$  columns where  $p$  is a small oversampling parameter. Let  $S^r = AR$  and  $S^c = A^*R$  be samples for the row (superscript  $r$ ) and column bases (superscript  $c$ ) of  $A$ , respectively. Algorithm 1 with  $R^r \equiv R^c \equiv R$  computes the HSS representation of  $A$  using the information available in the samples  $S^r$  and  $S^c$  by hierarchically compressing (using interpolative decompositions; see below) the off-diagonal blocks of  $A$ , starting from the leaves.

Let  $D_\tau$  for a nonleaf node  $\tau$  with children  $\nu_1$  and  $\nu_2$  be defined as

$$(2.5) \quad D_\tau = \begin{bmatrix} D_{\nu_1} & A_{\nu_1, \nu_2} \\ A_{\nu_2, \nu_1} & D_{\nu_2} \end{bmatrix}.$$

If  $\{\tau_1, \tau_2, \dots, \tau_q\}$  are all the nodes on level  $\ell$  of the HSS tree, then

$$(2.6) \quad D^{(\ell)} = \text{diag}(D_{\tau_1}, D_{\tau_2}, \dots, D_{\tau_q})$$

is an  $N \times N$  block diagonal matrix. The main idea of the randomized sampling Algorithm 1 is to construct a sample matrix  $S^{(\ell)}$  for each level of the tree as

$$(2.7) \quad S^{(\ell)} = (A - D^{(\ell)})R = S^r - D^{(\ell)}R.$$

This sample matrix  $S^{(\ell)}$  captures the action of a product of the block off-diagonal part of  $A$  with a set of random vectors  $R$ . It is exactly this block off-diagonal part that needs to be compressed using low-rank approximation.

Another crucial component of the randomized sampling algorithm is the interpolative decomposition (ID) [19]. The ID computes a factorization of a rank- $k$  matrix  $Y \in \mathbb{C}^{m \times n}$  by expressing  $Y$  as a linear combination of a set  $J$  of  $k$  selected columns of  $Y$ :

$$(2.8) \quad [X, J] = \text{ID}(Y), \quad \text{such that (s.t.)} \quad Y = Y(:, J)X, \quad Y(:, J) \in \mathbb{C}^{m \times k}, \quad X \in \mathbb{C}^{k \times n},$$

or it can be modified to take a compression tolerance  $\varepsilon$ , such that

$$(2.9) \quad [X, J] = \text{ID}(Y, \varepsilon), \quad \text{s.t.} \quad Y = Y(:, J)X + E, \quad Y(:, J) \in \mathbb{C}^{m \times k'}, \quad X \in \mathbb{C}^{k' \times n},$$

with  $\|E\| = \mathcal{O}(\varepsilon)$  and  $k' \leq k$  the  $\varepsilon$ -numerical rank with respect to (w.r.t.) the chosen norm. The ID can be computed from a rank-revealing or column pivoted QR decomposition [16, 45]

$$(2.10) \quad Y\Pi = Q \begin{bmatrix} R_1 & R_2 \end{bmatrix},$$

where  $R_1$  is upper-triangular and  $\Pi$  is a permutation matrix, followed by a triangular solve such that

$$(2.11) \quad Y = (QR_1) ([I \quad R_1^{-1}R_2] \Pi^{-1}) \equiv Y(:, J)X.$$

A consequence of using the ID in Algorithm 1 is that  $B_{\nu_1, \nu_2} = A(I_{\nu_1}^r, I_{\nu_2}^c)$  is a submatrix of the original matrix  $A$ . Furthermore, it also leads to a special structure for the  $U_\tau$  and  $V_\tau$  generators:

$$(2.12) \quad U_\tau = \Pi_\tau^r \begin{bmatrix} I \\ E_\tau^r \end{bmatrix} \quad \text{and} \quad V_\tau = \Pi_\tau^c \begin{bmatrix} I \\ E_\tau^c \end{bmatrix},$$

referred to as interpolative bases, which can be exploited in the computations. Note that these interpolative bases are not orthonormal. Although creating orthonormal bases might slightly improve stability, the interpolative structure improves performance of the compression algorithm and the ULV decomposition; see section 2.3.

**2.2.1. Adaptive scheme to determine the HSS-rank.** In practice however, the HSS-rank of the matrix is not known in advance. In this case, Algorithm 1 can be called repeatedly while increasing the number of columns of  $R$ ,  $S^r$ , and  $S^c$ . As long as  $d < r + p$ , the ID in line 9 will fail. Suppose the ID fails at node  $\tau$ , i.e., the required accuracy  $\varepsilon$  is not reached, but the descendants of node  $\tau$  are successfully compressed. In that case, during the next iteration of Algorithm 1 with  $d \leftarrow d + \Delta d$ , it is not necessary to redo the compression (ID) or the extraction of  $D$  and  $B$  for the descendants of node  $\tau$ . However, those descendants do have to update the  $\Delta d$  new columns in  $R^{r/c}$  (lines 12 and 15) and  $S^{r/c}$  (lines 5, 8 and 10). In [47], this adaptive rank scheme is presented in more detail.

**2.2.2. Implementation issues.** The random matrices  $R^r$  and  $R^c$  are filled element by element using a pseudorandom number generator. Our implementation offers the `minstd_rand` and `mt19937` generators from the C++11 standard while the distribution can be either uniform over  $[0, 1)$  or standard normal (Gaussian)  $\mathcal{N}(0, 1)$ .

---

**Algorithm 1:** Computing the HSS factorization of a nonsymmetric matrix.

---

```

1 Function  $A_{\text{hss}} = \text{HSSCompress}(R^r, R^c, S^r, S^c, \varepsilon, \tau = \text{root}(A_{\text{hss}}))$ 
   Data:  $S^r = AR^r$  and  $S^c = A^*R^c$  with  $\{S^r, S^c, R^r, R^c\} \in \mathbb{R}^{N \times d}$ ,
        $d \geq r_{\max} + p$ 
   Result:  $A_{\text{hss}}$ :  $D_\tau$  (leaves),  $B_{\nu_1, \nu_2}$ ,  $B_{\nu_2, \nu_1}$  (nonleaves),  $U_\tau$ ,  $V_\tau$  (all except
       root).

2 foreach  $\nu \in \text{child}(\tau)$  do  $\text{HSSCompress}(R^r, R^c, S^r, S^c, \nu)$ 
3 if  $\text{child}(\tau) \equiv \emptyset$  then
4    $D_\tau = A(I_\tau, I_\tau)$ 
5    $S_\tau^r = S^r(I_\tau, :) - D_\tau R^r(I_\tau, :)$             $S_\tau^c = S^c(I_\tau, :) - D_\tau^* R^c(I_\tau, :)$ 
6 else //  $\nu_1$  and  $\nu_2$  are the children of node  $\tau$ 
7    $B_{\nu_1, \nu_2} = A(I_{\nu_1}^r, I_{\nu_2}^c)$             $B_{\nu_2, \nu_1} = A(I_{\nu_2}^r, I_{\nu_1}^c)$ 
8    $S_\tau^r = \begin{bmatrix} S_{\nu_1}^r - B_{\nu_1, \nu_2} R_{\nu_2}^r \\ S_{\nu_2}^r - B_{\nu_2, \nu_1} R_{\nu_1}^r \end{bmatrix}$             $S_\tau^c = \begin{bmatrix} S_{\nu_1}^c - B_{\nu_2, \nu_1}^* R_{\nu_2}^c \\ S_{\nu_2}^c - B_{\nu_1, \nu_2}^* R_{\nu_1}^c \end{bmatrix}$ 
9    $[U_\tau^*, J_\tau^r] = \text{ID}((S_\tau^r)^*, \varepsilon)$             $[V_\tau^*, J_\tau^c] = \text{ID}((S_\tau^c)^*, \varepsilon)$ 
10   $S_\tau^r \leftarrow S_\tau^r(J_\tau^r, :)$             $S_\tau^c \leftarrow S_\tau^c(J_\tau^c, :)$ 
11 if  $\text{child}(\tau) \equiv \emptyset$  then
12    $R_\tau^r = V_\tau^* R^r(I_\tau, :)$             $R_\tau^c = U_\tau^* R^c(I_\tau, :)$ 
13    $I_\tau^r = I_\tau(J_\tau^r)$             $I_\tau^c = I_\tau(J_\tau^c)$ 
14 else
15    $R_\tau^r = V_\tau^* \begin{bmatrix} R_{\nu_1}^r \\ R_{\nu_2}^r \end{bmatrix}$             $R_\tau^c = U_\tau^* \begin{bmatrix} R_{\nu_1}^c \\ R_{\nu_2}^c \end{bmatrix}$ 
16    $I_\tau^r = [I_{\nu_1}^r \ I_{\nu_2}^r](J_\tau^r)$             $I_\tau^c = [I_{\nu_1}^c \ I_{\nu_2}^c](J_\tau^c)$ 

```

---

By default the linear congruential engine<sup>3</sup> `minstd_rand` is selected in combination with the Gaussian distribution.

The rank-revealing QR factorization, used in the ID, could be replaced by a *strong* rank-revealing QR factorization [28], with possibly greater accuracy and smaller HSS-rank but greater computational cost ( $\mathcal{O}(N^3)$ ). Note that the rank-revealing QR is applied to a matrix of reduced size, i.e.,  $\mathcal{O}(r \times r)$ , due to random sampling, so this additional computational cost might be negligible. This is left as future work. Two interesting alternative approaches to the randomized compression routine discussed in this section should be mentioned, namely adaptive cross approximation [10] and a matrix-free approach presented in [37].

**2.3. ULV-like factorization and solve.** Solving a linear system with an HSS matrix can be done by first computing a so-called ULV decomposition [18], where  $U$  and  $V^*$  are unitary matrices and  $L$  is lower triangular. However, in [55] and [57], the ULV decomposition is modified to take advantage of the special structure of the  $U_\tau$  and  $V_\tau$  generators; see (2.12). The resulting algorithm is referred to as ULV-like since it is no longer based on unitary transformations.

In the first step of a ULV factorization, zeros are introduced in the HSS block rows. This step can be done using, for instance, a full QL factorization

$$(2.13) \quad U_i = \Omega_\tau \begin{bmatrix} 0 \\ \tilde{U}_i \end{bmatrix}, \quad \Omega_\tau^* U_i = \begin{bmatrix} 0 \\ \tilde{U}_i \end{bmatrix}.$$

---

<sup>3</sup>This choice is motivated further in section 4.

However, thanks to the special structure of  $U_\tau$ , a multiplication from the left with a carefully chosen  $\Omega_\tau$  is much cheaper and has a similar effect:

$$(2.14) \quad \Omega_\tau = \begin{bmatrix} -E_\tau^r & I \\ I & 0 \end{bmatrix} \Pi_\tau^{rT} \rightarrow \Omega_\tau U_\tau = \Omega_\tau \Pi_\tau^r \begin{bmatrix} I \\ E_\tau^r \end{bmatrix} = \begin{bmatrix} 0 \\ I \end{bmatrix}.$$

We refer the reader to [47] for a detailed description of the ULV factorization and the corresponding solve.

**3. Multifrontal sparse LU factorization.** Our next object is to exploit the HSS matrix algebra in sparse factorizations. Most modern sparse factorization codes use either a supernodal or a multifrontal method. In both methods, panel factorization is performed on a sequence of dense submatrices, which, for example, correspond to the separators from a nested dissection ordering. After each panel factorization, a local Schur complement is formed. In the case of a supernodal algorithm, the local Schur complement is immediately scattered into the global Schur complement, whereas for a multifrontal method, the local Schur complement is stored and carried along temporarily, and its scattering to global Schur complement is delayed until that part of panel factorization is about to start. In view of the elimination tree capturing dataflow, a supernodal method transfers parts of the local Schur complement from the current node to a subset of ancestral nodes, whereas a multifrontal method requires data transfer only between the current node and the parent node. We believe it is easier to introduce HSS operations into a multifrontal method due to its simpler data transfer pattern.

This section briefly recalls the main ingredients of the multifrontal method for the LU factorization of general invertible sparse matrices. For a more detailed discussion of multifrontal methods, see [23, 38]. The method casts the factorization of a sparse matrix into a series of partial factorizations of many smaller dense matrices and Schur complement updates.

**3.1. Matrix reordering.** As a preprocessing step,  $A$  is first scaled and permuted for numerical stability:  $A \leftarrow D_r A D_c Q_c$ , where  $D_r$  and  $D_c$  are diagonal matrices that scale the rows and columns of  $A$  and  $Q_c$  is a column permutation that places large entries on the diagonal. We use the MC64 code by Duff and Koster [22] to perform the scaling and column permutation. Popular alternative scaling algorithms can be found in [48, 7, 20]. After that, a fill-reducing permutation  $A \leftarrow P A P^T$  is applied in order to reduce the number of nonzero elements in the LU factors. Permutation matrix  $P$  is computed using nested dissection applied to the adjacency graph of  $A + A^T$ , using one of the graph partitioning tools SCOTCH [44] or METIS [32]. Instead of nested dissection, other heuristics like AMD [5] can be used.

The multifrontal method relies on a structure called the *elimination tree*. The elimination tree serves as a task and data-dependency graph for both the factorization and the solution process. A few equivalent definitions of the elimination tree are available. We use the following, and we recommend the survey by Liu [38] for more detail on the method and the survey by L'Excellent for more detail about implementation issues like parallelism, memory usage, numerical aspects, etc. [36].

**DEFINITION 3.1.** Assume  $A = LU$ , where  $A$  is an  $N \times N$  sparse, structurally symmetric matrix. The elimination tree of  $A$  is a tree with  $N$  nodes, where the  $i$ th node corresponds to the  $i$ th column of  $L$  and with the parent relations defined by  $\text{parent}(j) = \min\{i : i > j \text{ and } \ell_{ij} \neq 0\}$  for  $j = 1, \dots, N-1$ .

In practice, nodes are amalgamated: nodes that represent columns and rows of

the factors with similar structures are grouped together in a single node. For instance, when using nested dissection reordering, all vertices from the same graph separator can be grouped in one elimination tree node. In the end, each node corresponds to a square dense matrix, referred to as a *frontal matrix*, with the following  $2 \times 2$  block structure:

$$(3.1) \quad \mathcal{F}_i = \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}.$$

**3.2. Numerical factorization.** Multifrontal factorization of the matrix consists of a bottom-up traversal of the tree, following a topological order (a node is processed before its parent). Processing a node means first forming (or *assembling*) the frontal matrix followed by elimination of the fully summed variables in the  $F_{11}$  block and finally a Schur complement update step. The frontal matrix  $\mathcal{F}_i$  is formed by summing the rows and columns of  $A$  corresponding to the variables in the  $F_{11}$ ,  $F_{21}$ , and  $F_{12}$  blocks, with the temporary data—the extended update matrices  $\bar{\mathcal{U}}_\nu$ —that have been produced by the children of  $i$  after their elimination step, i.e.,

$$(3.2) \quad \mathcal{F}_i = A_i + \sum_{\nu \in \text{child}(i)} \bar{\mathcal{U}}_\nu = \begin{bmatrix} A(I_i^{\text{sep}}, I_i^{\text{sep}}) & A(I_i^{\text{sep}}, I_i^{\text{upd}}) \\ A(I_i^{\text{upd}}, I_i^{\text{sep}}) & 0 \end{bmatrix} + \bar{\mathcal{U}}_{\nu_1} + \bar{\mathcal{U}}_{\nu_2} + \cdots,$$

where  $I_i = \{I_i^{\text{sep}}, I_i^{\text{upd}}\}$  is the set of row and column indices of  $\mathcal{F}_i$  w.r.t. the global matrix  $A$ , after reordering. Eliminating the fully summed variables in the  $F_{11}$  block is done through a partial factorization of  $\mathcal{F}_i$ , typically via a standard dense matrix factorization of the  $F_{11}$  block. Next, the Schur complement (contribution block or update matrix) is computed as  $\mathcal{U}_i = F_{22} - F_{21}F_{11}^{-1}F_{12}$  and stored in temporary memory. In contrast to the elimination step which uses straightforward dense matrix operations (high performance LAPACK/BLAS3 codes), the assembly step (3.2) requires index manipulation and indirect addressing while summing up  $\mathcal{U}_k$ . For example, if two children's update matrices  $\mathcal{U}_k = \begin{bmatrix} a_k & b_k \\ c_k & d_k \end{bmatrix}$ ,  $k = \nu_1, \nu_2$ , have subscript sets  $I_1^{\text{upd}} = \{1, 2\}$  and  $I_2^{\text{upd}} = \{1, 3\}$ , respectively, then those update matrices can only be added after aligning the index sets of the two matrices by padding with zero entries:

$$(3.3) \quad \mathcal{U}_1 \overset{\uparrow}{\leftarrow} \mathcal{U}_2 = \bar{\mathcal{U}}_1 + \bar{\mathcal{U}}_2 = \begin{bmatrix} a_1 & b_1 & 0 \\ c_1 & d_1 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} a_2 & 0 & b_2 \\ 0 & 0 & 0 \\ c_2 & 0 & d_2 \end{bmatrix} = \begin{bmatrix} a_1 + a_2 & b_1 & b_2 \\ c_1 & d_1 & 0 \\ c_2 & 0 & d_2 \end{bmatrix}.$$

This summation operation is called *extend-add*, denoted by  $\overset{\uparrow}{\leftarrow}$ . The relationship between frontal matrices and update matrices can be revealed by  $\mathcal{F}_i = A_i \overset{\uparrow}{\leftarrow} \mathcal{U}_{\nu_1} \overset{\uparrow}{\leftarrow} \mathcal{U}_{\nu_2} \overset{\uparrow}{\leftarrow} \cdots \overset{\uparrow}{\leftarrow} \mathcal{U}_{\nu_q}$ , where nodes  $\nu_1, \nu_2, \dots, \nu_q$  are the children of  $i$ .

Each partial factorization might involve pivoting within the frontal matrix. It can also happen that no suitable pivot can be found during a step of partial factorization. In this situation, the corresponding row and column remain unfactored and are sent to the parent node. This strategy is used, for instance, in the MUMPS [6] code. Currently, our code does not perform any such delayed pivoting, but instead relies on static pivoting (using MC64) and partial pivoting during the LU decomposition of the  $F_{11}$  blocks.

**3.3. Solution.** Once the factors are computed, the solution  $x$  of  $Ax = b$  is computed in two steps: forward solution by doing a triangular solution with the  $L$  factor and backward substitution by doing a triangular solution with the  $U$  factor. The forward solution step is a bottom-up topological traversal of the elimination tree, while the backward substitution is a top-down traversal.

**4. Multifrontal solver with HSS frontal matrices.** This section explains how a multifrontal solver (see section 3) can be used in combination with the HSS data-structures and algorithms from section 2 to improve the computational complexity and storage requirements. This section closely follows [55].

**4.1. Selection of HSS frontal matrices.** Note that the largest frontal matrices, those that determine the computational complexity of the solver, typically correspond to nodes closer to the root of the elimination tree. Let the top of the tree, i.e., the root node, be at level  $\ell = 0$  of the tree. Then, define a switch-level  $\ell_s$  such that the frontal matrices at levels  $\ell \geq \ell_s$  of the elimination tree are stored as regular dense matrices whereas those at levels  $\ell < \ell_s$  are compressed using the HSS format. According to the analysis in [55],  $\ell_s$  should be chosen such that the factorization costs above and below the switch-level are equal. However, this rule is not very practical, and experiments show that performance depends crucially on the choice of  $\ell_s$ .

**4.2. Separator reordering.** Apart from the scaling and permutation of  $A$  for stability, and nested dissection reordering to reduce fill-in, an additional reordering is applied to the index set of each separator. This reordering is needed to obtain favorable HSS rank structure in the corresponding frontal matrices. It is computed by recursively bisecting the graph of the separator into subgraphs of size approximately  $s$  (defaults to  $s = 128$ ), using a graph partitioning tool (SCOTCH or METIS). Each partition then corresponds to a leaf in the HSS tree of the corresponding frontal matrix. However, since a separator graph can be disconnected, it is enriched with length-two connections from the connectivity graph before it is passed to the partitioner; see also the discussion in [42]. Note that other reorderings can be used instead of nested dissection. The influence of the reordering on the ranks of off-diagonal blocks is studied in [53].

**4.3. Skinny extend-add.** From here on, we assume that a binary elimination tree is used. The steps followed for each HSS frontal matrix  $\mathcal{F}_i$  are as follows. First, a random matrix  $R_i \in \mathbb{C}^{\#I_i \times d_i}$  is constructed. If the children  $\nu_1$  and  $\nu_2$  of  $i$  are also HSS, then  $R_i$  is constructed as follows:

$$(4.1) \quad R_i(r, c) = \begin{cases} R_{\nu_1}(r, c) \equiv R_{\nu_2}(r, c) & \text{if } c < \min(d_{\nu_1}, d_{\nu_2}), I_i(r) \in I_{\nu_1}^{\text{upd}}, \text{ and } I_i(r) \in I_{\nu_2}^{\text{upd}}, \\ R_{\nu_1}(r, c) & \text{if } c < d_{\nu_1} \text{ and } I_i(r) \in I_{\nu_1}^{\text{upd}}, \\ R_{\nu_2}(r, c) & \text{if } c < d_{\nu_2} \text{ and } I_i(r) \in I_{\nu_2}^{\text{upd}}, \\ \text{random}(r, c) & \text{otherwise.} \end{cases}$$

The random matrices of the children are merged in the parent  $R_i$ , and any elements not present in any of the children's  $R$  are generated. This extend-merge procedure is illustrated in Figure 2. If node  $i$  has no (HSS) children,  $R_i$  is generated. However, it is important that corresponding “random” entries in  $R_{\nu_1}$  and  $R_{\nu_2}$  are equal, since that allows efficient evaluation of  $S_i^r = \mathcal{F}_i R_i$  (similarly for  $\mathcal{F}_i^* R_i$ ) based on

$$(4.2) \quad \mathcal{F}_i R_i = (A_i \overset{\uparrow}{\leftarrow} \mathcal{U}_{\nu_1} \overset{\uparrow}{\leftarrow} \mathcal{U}_{\nu_2}) R_i = (A_i R_i) \overset{\uparrow}{\leftarrow} (\mathcal{U}_{\nu_1} R_i(I_{\nu_1}^{\text{upd}}, :)) \overset{\uparrow}{\leftarrow} (\mathcal{U}_{\nu_2} R_i(I_{\nu_2}^{\text{upd}}, :)),$$

where  $R_i(I_{\nu_1}^{\text{upd}}, :)$  denotes the subset of rows of  $R_i$  which are also in  $I_{\nu_1}^{\text{upd}}$  and  $\overset{\uparrow}{\leftarrow}$  denotes an extend-add operation where the extend is only done for the rows, not the columns. By the construction of  $R_i$  (4.1), the first  $d_{\nu_1}$  columns of  $R_i(I_{\nu_1}^{\text{upd}}, :)$  are already available at node  $\nu_1$ , which is convenient for the evaluation of  $\mathcal{U}_{\nu_1} R_i(I_{\nu_1}^{\text{upd}}, :)$ . Evaluation of  $\mathcal{U}_{\nu_1} R_i(I_{\nu_1}^{\text{upd}}, :)$  is discussed in more detail in section 4.4. When generating

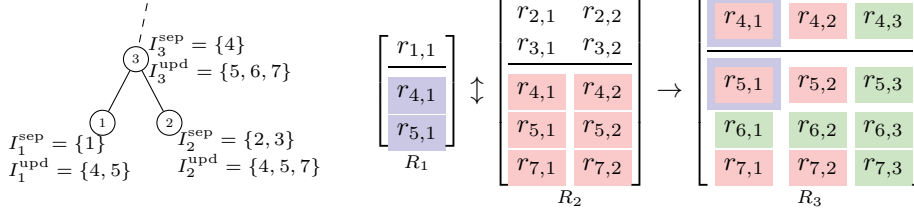


FIG. 2. Illustration of the extend-merge procedure for the random vectors. Node 3 needs three random vectors. It can get elements  $r_{4,1}$  and  $r_{5,1}$  from either child 1 or 2. Elements  $r_{7,1}$ ,  $r_{7,2}$ ,  $r_{4,2}$ , and  $r_{5,2}$  are copied from child 2. Elements  $r_{6,1}$  and  $r_{6,2}$  in  $R_3$  are generated with the properly seeded pseudo-random number generator. When the adaptive HSS compression scheme decides that a third column has to be added to  $R_3$ , those elements are also generated.

rows in  $R_i$ , the random number generator is seeded for each row using the global row index  $I_i(r)$  to ensure that  $R_i$  is consistent with its sibling. This frequent seeding is the reason the linear congruential pseudo-random engine `minstd_rand` was chosen as default over, for instance, the `mt19937` Mersenne-Twister, which has a much bigger internal state.

The frontal matrices  $\mathcal{F}_i$  with  $\text{level}(i) < \ell_s$  are completely approximated by HSS and are never explicitly formed as a dense matrix. This is in contrast to earlier, so-called partially structured approaches where, for instance, only the  $F_{11}$  or the  $F_{21}$ ,  $F_{11}$ , and  $F_{12}$  blocks are compressed [52]. Partially structured approaches typically at one point or another form a dense representation of the  $F_{22}$  block, perform the Schur complement update on it, and then use this dense update matrix in the extend-add procedure. This is done to avoid having to perform an overly complicated extend-add operation on HSS matrices. However, the approach followed here does not require first assembling a dense frontal matrix before doing HSS compression. This is due to the use of the randomized HSS compression, Algorithm 1, which only requires matrix-vector multiplication and extraction of selected elements from the frontal matrix.

When  $R_i$ ,  $S_i^r$ , and  $S_i^c$  have been constructed, HSS compression using Algorithm 1 can be performed. However, when  $d_i - p$  is less than the HSS-rank of  $\mathcal{F}_i$ , Algorithm 1 will fail. In that case, columns are added to  $R_i$ , i.e.,  $d_i \leftarrow d_i + \Delta d$  ( $\Delta d = 128$  by default), the new columns of  $S_i^r$  and  $S_i^c$  are computed, and Algorithm 1 is called again, this time only updating the new columns of  $R_i$ ,  $S_i^r$ , and  $S_i^c$ . Due to the use of the ID in Algorithm 1, HSS generators  $D_\tau$  and  $B_{\nu_1, \nu_2}$  are submatrices of  $\mathcal{F}_i$ . Hence, a routine to extract specific elements from  $\mathcal{F}_i$  is required. This routine will be described in section 4.5.

**4.4. ULV factorization and low-rank Schur complement update.** After HSS compression, a factorization of  $F_{i_{11}}$  is performed: classical row-pivoted LU if  $\mathcal{F}_i$  is dense, and ULV if it is HSS. For a dense frontal matrix,  $\mathcal{U}_i = F_{i_{22}} - F_{i_{21}} F_{i_{11}}^{-1} F_{i_{12}}$  is computed explicitly. In the HSS case,  $F_{i_{22}}$  is kept in HSS form and the update  $F_{i_{21}} F_{i_{11}}^{-1} F_{i_{12}} = \Theta_i^* \Phi_i$  is stored as a low-rank product. Expressions for  $\Theta_i^*$  and  $\Phi_i$  are derived and presented in detail in [55] for symmetric and in [57] for nonsymmetric matrices. Given  $\mathcal{U}_i = F_{i_{22}} - \Theta_i^* \Phi_i$ , the multiplication with  $\mathcal{U}_i$  in (4.2) can be performed efficiently using HSS matrix-vector multiplication for  $F_{i_{22}}$  and two dense (rectangular) matrix products for  $\Theta_i^*$  and  $\Phi_i$ .

**4.5. Extracting elements from an HSS matrix.** Finally, extracting elements from  $\mathcal{F}_i$  requires extracting elements from an HSS matrix. In [55] a routine is pre-

sented for extracting multiple elements from an HSS matrix while trying to minimize the number of traversals through the HSS tree. We use a conceptually simpler algorithm based on the HSS matrix-vector multiplication. By multiplying an HSS matrix with unit vectors, selected columns can be extracted. At the leaf nodes, instead of multiplying with a unit vector, one can simply select the proper columns of  $V^*$ . Unlike for matrix-vector multiplication, during element extraction parts of the tree traversal can be pruned.

**4.6. Preconditioning versus iterative refinement.** Direct solvers often use a few steps of iterative refinement to improve the solution quality [54]. However, the multifrontal method with HSS compression as presented in this paper is used as a preconditioner for GMRES instead. For the same number of multifrontal solve steps (preconditioner applications), a Krylov solver typically leads to smaller residuals than iterative refinement. This is particularly useful when the HSS compression tolerance is increased, since in that case the quality of the HSS-multifrontal preconditioner decreases and the number of outer iterations increases. Iterative refinement might not converge in this case since it no longer defines a contraction.

**4.7. Solver complexity.** The computational complexity of a standard multifrontal solver is typically dominated by the dense linear algebra corresponding to the few largest frontal matrices. For instance, a nested dissection reordering on a  $d$ -dimensional mesh with  $N = k^d$  vertices has a top separator with  $\mathcal{O}(k^{d-1})$  vertices, leading to an overall complexity of  $\mathcal{O}(k^{3(d-1)})$ , i.e.,  $\mathcal{O}(N^{3/2})$  and  $\mathcal{O}(N^2)$  for two-dimensional (2D) and three-dimensional (3D) meshes, respectively.

For the HSS-embedded multifrontal solver, the complexity is dominated by the HSS compression of the dense frontal matrices, which in turn depends on the rank pattern. Earlier works by Chandrasekaran et al. [17] and Engquist and Ying [25] showed the rank patterns of the elliptic and the Helmholtz operators, respectively. Xia showed complexities for the randomized HSS multifrontal solver assuming different rank patterns [55]. Combining the above results, we summarize the solver complexities for two types of PDEs and two sparse solvers in Table 1. A major bottleneck for direct solvers is often the large memory usage requirement. The HSS-embedded multifrontal solver has lower asymptotic memory usage than the traditional multifrontal solver, as also illustrated in Table 1. For the two PDE problems considered in Table 1 the HSS-embedded solver has optimal memory scaling.

TABLE 1

*Summary of the complexities of the standard multifrontal solver (MF) and the randomized HSS-embedded multifrontal solver (MF-HSS-RS) applied to two important classes of problems. The mesh size per side is  $k$  and the matrix dimensions are  $N = k^2$  in 2D and  $N = k^3$  in 3D.*

Problem		HSS rank	MF		MF-HSS-RS	
			Factor flops	Memory	Factor flops	Memory
2D ( $k \times k$ )	elliptic	$\mathcal{O}(1)$	$\mathcal{O}(N^{3/2})$	$\mathcal{O}(N \log N)$	$\mathcal{O}(N)$	$\mathcal{O}(N)$
	Helmholtz	$\mathcal{O}(\log k)$				
3D ( $k \times k \times k$ )	elliptic	$\mathcal{O}(k)$	$\mathcal{O}(N^2)$	$\mathcal{O}(N^{4/3})$	$\mathcal{O}(N \log N)$	$\mathcal{O}(N)$
	Helmholtz	$\mathcal{O}(k)$				

**5. Shared memory parallel implementation.** The algorithm presented in section 4 has been implemented using C++ and OpenMP, targeting shared memory platforms. The code relies on BLAS, LAPACK, METIS, and/or SCOTCH and a recent C++11 compliant compiler with support for OpenMP 3.1 or higher. The code

## LISTING 1

*Bottom-up topological parallel tree traversal implemented with recursion and the OpenMP (3.1) task construct.*

```

void Tree::postorder(depth=0) {
    if (depth < d_max) {
        if (left_child)
#pragma omp task untied default(shared) final(depth >= d_max-1) mergeable
            left_child->postorder(depth+1)
        if (right_child)
#pragma omp task untied default(shared) final(depth >= d_max-1) mergeable
            right_child->postorder(depth+1)
#pragma omp taskwait
    } else {
        if (left_child) left_child->postorder(depth+1)
        if (right_child) right_child->postorder(depth+1)
    }
    do_stuff(depth); // factor/compress..., can generate more tasks
}

```

makes heavy use of the OpenMP task construct. OpenMP was chosen because it is easy to use, performs well, and is well documented and supported. However, alternatives like Intel Threading Building Blocks [46] or Cilk(+) [11] offer conceptually similar task parallelism. Switching to one of those should not be hard. While other runtime systems like QUARK [58], DAGuE/ParSEC [14], and StarPU [9] (distributed memory task scheduling) and OmpSs [24] might have certain specific advantages over the OpenMP runtime, many of those innovations (for instance, explicit modeling of task dependencies or task-offloading) are eventually incorporated in the OpenMP standard as well.

OpenMP tasks are created and scheduled at runtime by the scheduler. Task schedulers typically use a work stealing [12] or task stealing strategy to balance load between threads. Each thread/core has its own local queue of tasks. When a thread runs out of work it can steal a task from one of the other thread's task queues.

**5.1. Task based tree parallelism.** Traversals of both the elimination tree and the HSS hierarchy allow for tree parallelism; i.e., independent subtrees can be processed concurrently. For instance, multifrontal factorization requires bottom-up topological traversal of the elimination tree, just like HSS compression requires bottom-up traversal of the HSS hierarchy. The code in Listing 1 shows how to do a parallel bottom-up tree traversal using the OpenMP task construct. The tree is stored as objects of a class `Tree` with two members `left_child` and `right_child`, both pointers to subtrees, also objects of type `Tree`. In Listing 1, the variable `depth` keeps track of the recursion depth; and no more tasks are generated after a certain depth to avoid excessive overhead of creating too fine-grained tasks. Experiments show that setting `d_max` to  $\log_2(\#\text{threads}) + 3$  leads to a good task granularity. With this setting, the maximum number of tasks at any given point in time is about  $2^{d_{\max}} = 8 \cdot \#\text{threads}$ . This is enough to ensure good load balance and avoids excessive task creation overhead. OpenMP tasks supports an `if` clause, so the check `if(depth < d_max)` could have been put in the OpenMP pragma. However, optimizing the code to perform this check outside the directive completely avoids all task creation and synchronization overhead when it evaluates to false. The `final(condition)` clause informs the OpenMP runtime that the generated task will not generate more tasks if `condition` evaluates to true. Finally, the `untied` clause informs the runtime that this task can be moved to a different thread when it encounters a scheduling point. For instance,

when a task spawns a new task, the spawning task may be moved to another thread. Untied tasks allow for better load balance, whereas tied tasks (the default) typically lead to better data locality. The `taskwait` pragma ensures that processing of the children is finished before continuing with the parent.

**5.2. Hybrid node and tree parallelism.** Exploiting tree parallelism alone as in Listing 1 does not scale well due to the limited degree of parallelism near the root. Although the HSS-multifrontal algorithm can exploit two nested levels of tree parallelism (elimination tree and HSS hierarchy), the scaling bottleneck remains. To overcome this, one needs to exploit parallelism in the computational work inside the tree nodes, which are mostly dense linear algebra operations. However, work sharing constructs like OpenMP `parallel for` loops are not allowed within OpenMP tasks. Moreover, calling multithreaded BLAS or LAPACK routines from multiple tasks/threads leads to oversubscription and generally poor performance. This is because existing multithreaded BLAS/LAPACK libraries are optimized to use the entire machine. One possible strategy is to exploit tree parallelism only for the lower levels of the tree and switch to a sequential processing of the nodes higher up in the tree while switching to multithreaded linear algebra. However, this leads to many synchronization points and does not scale with an increasing number of threads. Our approach, on the other hand, is to use task parallelism within the tree nodes as well to allow for a seamless transition between tree and node parallelism, since scheduling of tasks is left to the runtime system. When getting closer to the root node, there is a shift from tree to node parallelism. This is illustrated in Figure 3. Even in the case of highly unbalanced trees, the runtime can assign work evenly to the available cores. We chose not to use an existing library for the task based dense linear algebra, for instance, PLASMA (based on the QUARK runtime), since we wished to exploit the same threading mechanism (OpenMP) already used for the tree parallelism.

**5.3. Parallel BLAS and LAPACK.** One of the most time consuming operations of the algorithm is dense matrix-matrix multiplication  $C \leftarrow \alpha AB + \beta C$ . This can be implemented easily with recursion and task parallelism [41], by splitting the problem into smaller matrix-matrix multiplications; this strategy is referred to as divide-and-conquer and is often used in so-called cache-oblivious algorithms [26]. How the matrices are split depends on their shapes. Let  $A$  be  $m \times k$  and  $B$  be  $k \times n$ ; then

$$(5.1) \quad C \leftarrow \begin{cases} \alpha AB + \beta C & \text{if } m \times n \times k \leq T, \\ \alpha \begin{bmatrix} AB_0 & AB_1 \end{bmatrix} + \beta \begin{bmatrix} C_0 & C_1 \end{bmatrix} & \text{else if } n \geq \max(m, k), \\ \alpha \begin{bmatrix} A_0 B \\ A_1 B \end{bmatrix} + \beta \begin{bmatrix} C_0 \\ C_1 \end{bmatrix} & \text{else if } m \geq k, \\ \alpha (A_0 B_0 + A_1 B_1) + \beta C & \text{else.} \end{cases}$$

The last case in (5.1), short fat  $A$  times tall skinny  $B$ , uses two consecutive recursive matrix-matrix multiplication calls. Cases 2 and 3 start two multiplications in parallel, spawning two tasks. The recursion ends when reaching case 1, with  $T$  a tuning parameter set by default to  $T = 64^3$ , where a sequential vendor optimized BLAS3 `*gemm` routine is called. Depending on the scalar type, one of four inlined template specialization functions for `gemm<scalar>` is executed to pick the correct version: `sgemm`, `dgemm`, `cgemm`, or `zgemm`. For the other BLAS2/3 routines that are required (for instance, triangular matrix multiplication and solve), a similar recursive approach

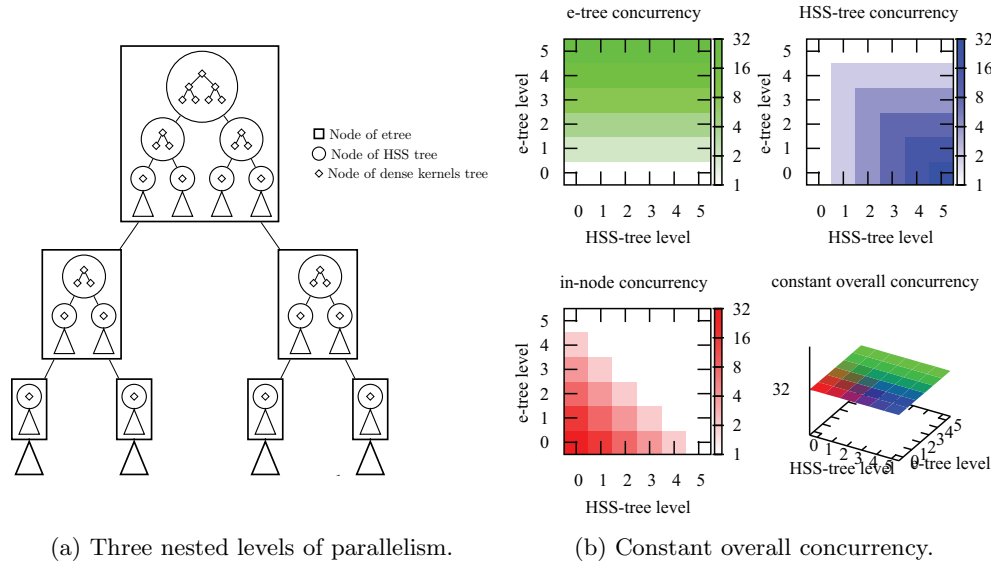


FIG. 3. Schematic illustration of the different types of concurrency in the code and the gradual shift from tree parallelism to in-node parallelism. (a) Tasks for dense kernels ( $\diamond$ ) are nested in nodes ( $\circ$ ) of the HSS trees, which are nested in the elimination tree nodes ( $\square$ ) (e-tree). (b) Left-to-right, top-to-bottom: (1) Elimination tree concurrency decreases when getting closer to the root node. (2) Closer to the root of the elimination tree, more HSS tree concurrency is exploited as it becomes available, i.e., while moving down the HSS tree away from the root. (3) Towards the root of the HSS tree and the root of the elimination tree, more in-node concurrency (parallel tasked dense algebra) is exploited. (4) The product of the three types of concurrency, i.e., the overall concurrency, remains constant throughout both the elimination and HSS trees.

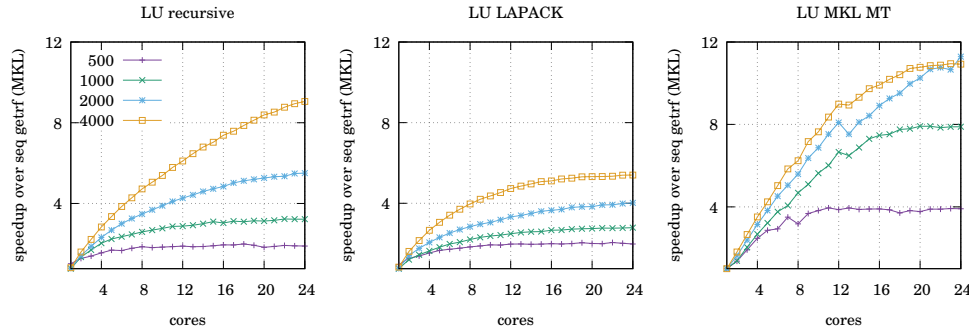


FIG. 4. Speedups over sequential `getrf` from Intel MKL for matrices of size  $500^2$  to  $4000^2$ . Left: Recursive LU decomposition using OpenMP tasked BLAS code. Middle: Reference LAPACK `getrf` using OpenMP tasked BLAS code. Right: MKL optimized multithreaded `getrf`. Our recursive implementation scales better than the reference netlib `getrf` with parallel BLAS but worse than the MKL optimized code. However, calling MKL multithreaded `getrf` from multiple threads simultaneously would lead to oversubscription and performance penalty. This is not a problem with the recursive LU because it uses the OpenMP task runtime, just like the rest of the code.

is used. This recursive task generation is also stopped when the recursion depth becomes too large, with the same `depth` parameter being passed through the entire code and incremented each time it enters a new task.

The code also requires some LAPACK functionality, namely LQ, LU, and RRQR decompositions. For those, we modify the reference Fortran LAPACK implementation to make use of our parallel (tasked) BLAS routines. Some vendor optimized LAPACK libraries not only use the LAPACK reference code on top of multithreaded BLAS calls, but also add additional optimizations to the LAPACK routines. Unfortunately, in our approach we cannot take advantage of these optimized multithreaded codes. Consider partial pivoted LU decomposition, used for the  $F_{11}$  block of a dense frontal matrix. Apart from the LAPACK `*getrf` routine using our OpenMP tasked BLAS routines, we also implemented a recursive LU factorization algorithm [21]. The parallelism in this algorithm has to come from the BLAS routines, triangular solve, row permutation, and matrix-matrix multiply. Figure 4 compares the performance and scalability of the two LU decomposition approaches with the MKL optimized implementation and shows that our implementation of LU scales nearly as well as MKL without sacrificing the ability to exploit subtree concurrency. A more scalable approach [15], based on so-called tiled algorithms instead of recursion, partitions the matrices into tiles of fixed sizes and assigns tasks to each of the tiles while explicitly modeling the data dependencies between the tasks. A DAG scheduler then executes the tasks while respecting their dependencies. OpenMP supports explicit task dependencies since version 4.0.<sup>4</sup> We intend to exploit this feature in the future to achieve more scalable dense linear algebra operations. For the rank-revealing QR decomposition we use a modified version of the LAPACK `*geqp3` code [45], a BLAS3 version of column pivoted QR. The routine is modified to call our parallel tasked BLAS and an extra tolerance parameter  $\varepsilon$  is added to stop the rank-revealing process as soon as the  $\varepsilon$ -rank has been found instead of computing the full decomposition. More precisely, numerical rank  $i$  is detected when  $R_{i+1,i+1}/R_{11} \leq \varepsilon$ , where  $R$  is the upper-triangular factor.

**5.4. Scaling bottlenecks.** Before the actual numerical factorization step, but after matrix scaling and nested dissection reordering, a symbolic factorization step is performed. During this step some memory is allocated, and the index sets  $I_\tau^{\text{upd}}$  are assembled. The symbolic factorization is a bottom-up tree traversal which is done in parallel, as in Listing 1. In a multithreaded setting, memory allocation can become a serious scaling bottleneck. We have found that the use of a scalable memory allocator, like TCMalloc [27] or the TBB scalable memory allocator [46] greatly improves the performance over, for instance, the default `malloc` in glibc.<sup>5</sup> For instance, running on a 60 core Intel Xeon Phi, the symbolic factorization phase runs up to  $56\times$  faster when using TBBMalloc instead of the default allocator.

**6. Numerical experiments.** This section presents various numerical results. Section 6.1 first focuses on some PDE problems on regular grids as this allows us to easily change the problem size. The following sections consider other matrices from various applications. Unless otherwise stated, the experiments are performed on a single 12-core socket of a single node of the NERSC Edison machine.<sup>6</sup> A compute node has two 12-core Intel Ivy Bridge processors at 2.4GHz. Double precision peak performance is 19.2Gflop/s per core, 230.4Gflop/s per socket, or 460.8Gflop/s per node. Each socket has 32GB DDR3 1866MHz memory, hence 64GB per node, with a STREAM [40] bandwidth of 48.5GB/s. We use the Intel 15.0.1 compiler with

<sup>4</sup>Not all compilers currently support the latest OpenMP 4.0 standard.

<sup>5</sup><http://www.gnu.org/software/libc/>

<sup>6</sup><https://www.nersc.gov/users/computational-systems/edison/>

sequential MKL.

**6.1. PDEs on a regular grid.** We start with a number of benchmarks for well-known PDEs on regular 2D and 3D grids to study scaling of time-to-solution, number of floating point operations, memory usage, HSS-ranks, etc., w.r.t. problem size. For these regular grids, a geometric nested dissection code is used instead of the default METIS graph partitioner. The following benchmark problems are considered:

- Poisson equation  $-\Delta u = f$  on a 2D grid (P2D) using the standard 5-point finite difference stencil with homogeneous Dirichlet boundary conditions.
- Poisson equation on a 3D grid (P3D) using the standard 7-point stencil with homogeneous Dirichlet boundary conditions.
- Convection diffusion equation [43]  $-\nu\Delta u + \mathbf{v} \cdot \nabla u = f$  on a 2D grid (C2D) using a 5-point upwind stencil, with viscosity  $\nu = 10^{-4}$  and

$$(6.1) \quad \mathbf{v} = (x(1-x)(2y-1) \quad y(1-y)(2x-1))^T.$$

- Convection diffusion, similar to the above, on 3D grid (C3D) with

$$(6.2) \quad \mathbf{v} = (2x(1-x)(2y-1)z \quad -y(1-y)(2x-1) \quad -(2x-1)(2y-1)z(1-z))^T.$$

- Helmholtz equation

$$(6.3) \quad (-\Delta - \omega^2/v(x)^2) u(x, \omega) = s(x, \omega)$$

on a 2D grid (H2D), with  $\omega$  the angular frequency,  $v(x)$  the seismic velocity, and  $u(x, \omega)$  the time-harmonic wavefield solution to the forcing term  $s(x, \omega)$ . The discretization uses a 9-point stencil and the frequency is set at  $f = 10\text{Hz}$  with  $\omega = 2\pi f$ . The seismic speed is  $v(x) = 1500\text{ m/s}$ . We use a sampling rate of about 15 points per wavelength and PML boundary conditions. This example is indefinite and uses complex arithmetic.

- Same as H2D, but 3D using a 27-point stencil (H3D).

A crucial parameter for performance is the number of levels  $\ell_s$  of the elimination tree for which HSS compression is performed. We call this the HSS switching level. Note that  $\ell_s = 0$  corresponds to a pure multifrontal solver. Unfortunately, the optimal  $\ell_s$  is impossible to predict a priori, so it is determined experimentally and will always be mentioned with each result. The same applies to the compression tolerance  $\varepsilon$ . In [55, Theorem 4.2], Xia suggests setting  $\ell_s$  such that the cost of factorization for the levels above the switch-level  $\ell_s$  equals the cost of factorization for the levels below the switch-level  $\ell_s$ . Unfortunately, since the HSS-ranks are not known a priori, this is not a practical guideline. Our recommendation is for the user to experiment with the values of  $\ell_s$  and  $\varepsilon$  to get some intuition as to appropriate values for the particular application. When  $\ell_s > 0$ , i.e., with HSS compression, the multifrontal solver is used as a preconditioner for restarted GMRES(30) with modified Gram-Schmidt and a zero initial guess. Without HSS compression, iterative refinement with the direct solver is used. All experiments are performed in double precision with relative or absolute stopping criteria  $\|u_i\|/\|u_0\| \leq 10^{-6}$  or  $\|u_i\| \leq 10^{-10}$ , where  $u_i = M^{-1}(Ax_i - b)$ , with  $M$  the approximate multifrontal factorization of  $A$ , is the preconditioned residual. The right-hand-side is always set to  $A \begin{bmatrix} 1 & 1 & \cdots & 1 \end{bmatrix}^T$ .

**6.1.1. Performance results.** Figure 5 shows timing results for the 2D (top) and 3D (bottom) Poisson equations on  $5000^2$  and  $125^3$  grids, respectively. Figure 5a shows numerical factorization time as a function of the number of levels in the elimination

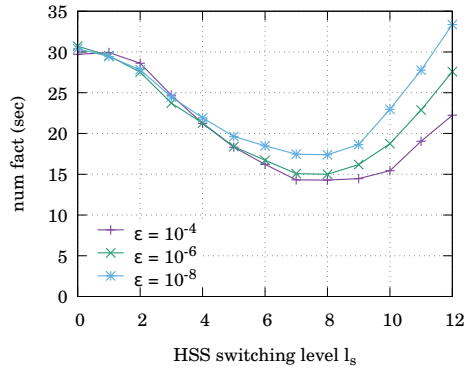
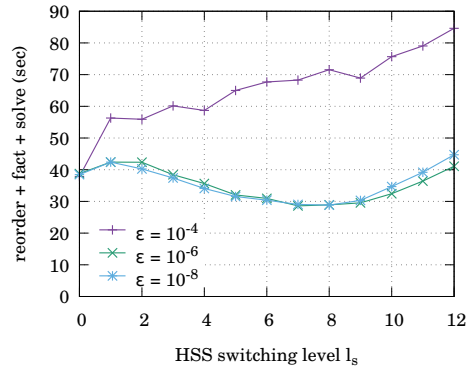
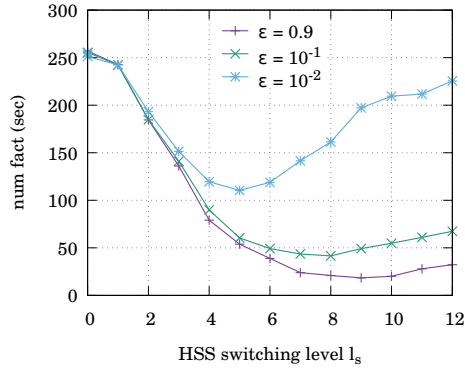
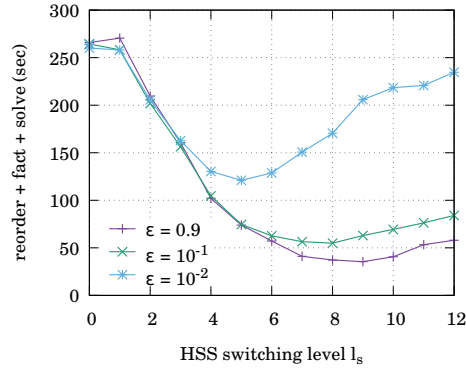
(a) Num fact time,  $5000^2$ .(b) Total solve time,  $5000^2$ .(c) Num fact time,  $125^3$ .(d) Total solve time,  $125^3$ .

FIG. 5. Times for factorization (left) and solve (right) of the 2D (top) and 3D (bottom) Poisson equations on  $5000^2$  and  $125^3$  grids as function of the number of levels  $\ell_s$  in the elimination tree for which HSS compression is applied. Different curves correspond to different HSS compression tolerances  $\varepsilon$ . For 3D, much more aggressive HSS compression can be used.

tree for which HSS compression was used. The HSS levels always correspond to the top levels of the elimination tree. This shows that applying HSS compression leads to a speedup of about  $2\times$  for 7 HSS levels. Different lines correspond to different HSS compression tolerances  $\varepsilon$ . Somewhat larger factorization speedups are possible for  $\varepsilon \geq 10^{-4}$ . However, this does not lead to faster time-to-solution. Figure 5b shows the cumulative time for nested dissection reordering, symbolic factorization, numerical factorization, and GMRES solve. For  $\varepsilon \geq 10^{-4}$ , the number of GMRES iterations, and thus the number of applications of the multifrontal solve, increases too much to get overall speedup; see also Table 2. Best results were obtained with  $\ell_s = 8$ ,  $\varepsilon = 10^{-7}$ , and only 2 GMRES iterations (3 multifrontal solves). Figures 5c and 5d show the timings for the 3D Poisson problem. For the 3D problem, much more aggressive HSS compression can be used. Best results were obtained with  $\ell_s = 10$ ,  $\varepsilon = 0.9$ , and 61 GMRES iterations. For the Poisson problem it seems that for 2D the direct solver is very efficient, with a modest speedup from HSS, while for 3D the HSS

TABLE 2  
GMRES iterations corresponding to the experiments shown in Figure 5.

$\varepsilon \setminus \ell_s$	0	1	2	3	4	5	6	7	8	9	10	11	12	
2D 5000 <sup>2</sup>	$10^{-4}$	1	8	9	13	14	18	19	22	23	23	24	24	25
	$10^{-6}$	1	2	3	3	3	3	3	3	3	3	3	3	3
	$10^{-8}$	1	2	2	2	2	2	2	2	2	2	2	2	2
3D 125 <sup>3</sup>	0.9	1	21	27	31	39	43	47	57	67	71	88	111	113
	$10^{-1}$	1	9	16	18	23	27	30	34	41	43	47	50	55
	$10^{-2}$	1	8	10	12	15	18	19	20	22	21	22	23	23

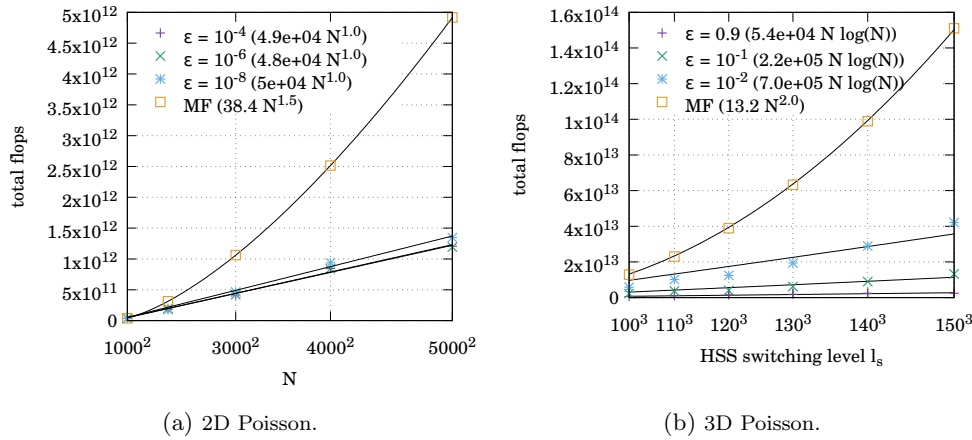


FIG. 6. Scaling of the number of floating point operations required to factor and solve a 2D(a) or 3D(b) Poisson equation. (a) The theory predicts  $\mathcal{O}(N^{3/2})$  complexity for the multifrontal (MF) solver and optimal  $\mathcal{O}(N)$  [55] complexity with HSS compression. (b)  $\mathcal{O}(N^2)$  complexity for the multifrontal solver and slightly lower complexity with HSS compression. The fits (black lines) are very sensitive to the data and not very reliable. However, note the smaller exponents and the larger constants for the new solver.

enabled factorization leads to a good preconditioner.

Figure 6a shows the total number of flops (numerical factorization and GMRES solve) for solving a 2D Poisson equation as function of the number of degrees of freedom, again for different compression tolerances. For the pure multifrontal method (no compression), the number of flops is  $\mathcal{O}(N^{3/2})$ , as predicted by the theory. For 2D Poisson, the HSS-rank is independent of the grid size [17], which leads to an optimal solver, i.e., linear scaling in the number of unknowns; see the fit in Figure 6a. Note the much larger constant for the HSS method. For the 5000<sup>2</sup> problem there is a reduction in the number of flops by a factor of about 3.3 $\times$ . However, the observed speedup (Figure 5b) is smaller than that. This is due to the fact that although the number of flops for the factorization decreases, the number of flops for the solution phase (and GMRES iterations) increases. Although multifrontal solve requires an order of magnitude less flops than factorization, it runs at much lower flop rates on modern hardware because it is limited by the memory bandwidth instead of the floating point unit. Additionally, the flop rate in the factorization phase is lower when using HSS compression due to the more fine-grained task decomposition. Figure 6b shows number of flops-to-solution for the 3D Poisson equation. For the very aggressive compression,  $\varepsilon = 0.9$ , the number of floating point operations for the 125<sup>3</sup> problem is reduced to 4.4% of the number of flops for the multifrontal method. Figure 7a

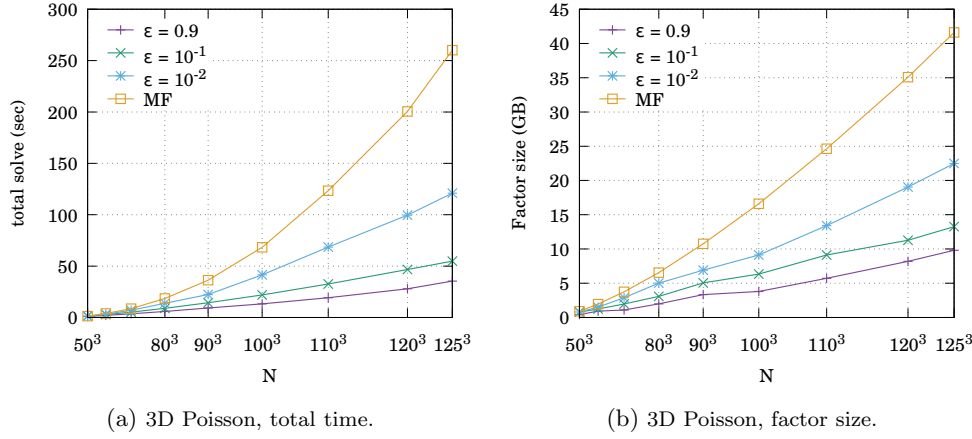


FIG. 7. (a) Total solve time (reordering, factorization, and solution) for the 3D  $125^3$  Poisson equation. (b) Factor size for the same problem. Different lines correspond to different HSS compression tolerance  $\epsilon$ , MF refers to pure multifrontal. The HSS enabled solver is faster for larger problems, and it allows us to solve larger problems.

TABLE 3

Number of levels  $\ell_s$  in the separator tree for which HSS compression is used, for the experiments shown in Figures 6 and 7. The number of levels  $\ell_s$  is chosen as the optimal.

$\epsilon \setminus N$	$1000^2$	$2000^2$	$3000^2$	$4000^2$	$5000^2$
$10^{-4}$	3	6	8	8	9
$10^{-6}$	3	6	7	8	9
$10^{-8}$	3	5	7	8	9

$\epsilon \setminus N$	$100^3$	$110^3$	$120^3$	$130^3$	$140^3$	$150^3$
0.9	9	9	10	11	11	11
$10^{-1}$	8	8	8	9	10	10
$10^{-2}$	5	6	6	6	7	7

shows the total solve time for the 3D problem for different grid sizes. These times include matrix reordering, factorization, and GMRES solve. Figure 7b shows the size of the factors. Table 3 shows  $\ell_s$ , the number of levels in the separator tree for which HSS compression is used, for the experiments shown in Figures 6 and 7. The optimal number of HSS levels slowly increases as the problem size increases.

Table 4 shows detailed results for the six PDE problems. The best speedups are obtained for the 3D problems. The code achieves good performance in flops per second for the factorization phase—although slightly less so for the HSS enabled code. Since the performance of the solve phase is not bounded by the floating point unit but rather by the memory bandwidth, we report the approximate attained bandwidth. The detailed results from Table 4 are summarized in Figure 8.

The  $\epsilon$  and  $\ell_s$  values used for Table 4 and Figures 5–7 were chosen to minimize the total time to factor and solve a single linear system, i.e., the optimal trade-off between factorization time and number of GMRES iterations. When multiple consecutive solves with the same matrix are required, one needs to select different  $\ell_s$  and  $\epsilon$  values. For many consecutive and highly accurate solves, the pure (exact) multifrontal factorization is probably optimal as it minimizes the number of multifrontal triangular solves. However, suppose only a few digits of accuracy are required. The

TABLE 4

Comparison of the standard multifrontal solver and the multifrontal solver with HSS compression for a number of PDEs on regular grids. All experiments are run on a 12-core Intel Ivy Bridge (peak 230.4Gflop/s and 48.5GB/s) in double precision. The code achieves good performance in terms of Gflop/s (for the factorization) or GByte/s (for the solve) and HSS compression leads to nice speedups over the standard multifrontal solver. A geometric nested dissection code is used for these regular grid problems.

Problem		P2D	P3D	C2D	C3D	H2D	H3D
Grid size		5000 <sup>2</sup>	125 <sup>3</sup>	5000 <sup>2</sup>	125 <sup>3</sup>	4000 <sup>2</sup>	100 <sup>3</sup>
Multifrontal	Nested dissection time (s)	2.5	0.25	2.1	0.24	2.9	0.43
	Symbolic factorization time (s)	3.6	8.0	3.4	8.1	4.5	6.0
	Factorization time (s)	29.1	254.7	28.8	252.4	53.6	259.1
	Factorization flops ( $\times 10^{12}$ )	4.9	50.0	4.9	50.0	10.1	53.5
	flop rate ( $\times 10^9$ Gflop/s)	168.4	196.3	170.1	198.1	188.4	206.5
	fraction of peak	73%	85%	74%	86%	82%	90%
	Factor size (GB)	28.4	41.6	28.4	41.6	36.3	35.5
	Solution time (s)	1.4	1.1	1.5	1.1	2.4	0.91
	Solution flops ( $\times 10^9$ )	7.9	10.5	7.9	10.5	21.1	18.2
	Solution bandwidth (GB/s)	20.3	37.8	18.9	37.8	15.1	39.0
	fraction of peak	42%	78%	39%	78%	31%	80%
	Total flops ( $\times 10^{12}$ )	4.9	50.0	4.9	50.0	10.1	53.5
	Total time (s)	36.6	264.1	35.8	261.8	63.4	266.4
Multifrontal + HSS	Nested dissection time (s)	2.5	0.26	2.1	0.24	2.9	0.43
	Separator reordering time (s)	1.1	0.40	1.1	0.35	1.3	0.68
	Symbolic factorization time (s)	2.0	0.55	2.1	0.8	2.9	1.8
	Factorization time (s)	14.5	19.6	13.6	41.5	30.5	92.8
	Factorization flops ( $\times 10^{12}$ )	1.5	2.0	1.3	5.0	3.9	18.0
	flop rate ( $\times 10^9$ Gflop/s)	103.4	102.0	95.6	120.5	127.9	194.0
	fraction of peak	45%	44%	41%	52%	56%	84%
	Factor size (GB)	22.5	9.8	21.6	14.6	29.6	21.4
	fraction of multifrontal	<b>79%</b>	<b>24%</b>	<b>76%</b>	<b>35%</b>	<b>82%</b>	<b>60%</b>
	Solution time (s)	4.1	15.3	4.3	75.9	22.5	71.2
	GMRES(30) iterations	3	67	3	234	11	152
	Solution flops ( $\times 10^9$ )	25.6	169.9	23.7	876.3	210.5	1,759.0
	Solution bandwidth (GB/s)	22.0	43.6	20.1	45.2	15.8	46.0
	fraction of peak	45%	90%	41%	93%	33%	95%
	HSS levels $\ell_s$ (total)	7 (22)	8 (18)	8 (22)	7 (18)	7 (22)	4 (18)
	HSS-rank	48	46	50	397	139	30
	HSS compression tolerance $\varepsilon$	$10^{-6}$	0.9	$10^{-5}$	0.1	$10^{-4}$	0.9
	Total flops ( $\times 10^{12}$ )	1.5	2.2	1.3	5.9	4.1	19.8
	fraction of multifrontal	<b>30%</b>	<b>4.4%</b>	<b>27%</b>	<b>12%</b>	<b>41%</b>	<b>37%</b>
	Total time (s)	24.2	36.1	23.2	118.8	60.1	166.9
	Speedup	<b>1.52</b> $\times$	<b>7.14</b> $\times$	<b>1.54</b> $\times$	<b>2.22</b> $\times$	<b>1.05</b> $\times$	<b>1.59</b> $\times$

multifrontal HSS solver can be used as a direct solver, and due to the smaller factor size the solve phase will be faster than a solve with the pure multifrontal code.

**6.1.2. Memory usage.** Figure 7b, showing the factor size for the 3D Poisson problem, illustrates the benefit of lower memory usage for the HSS enabled solver

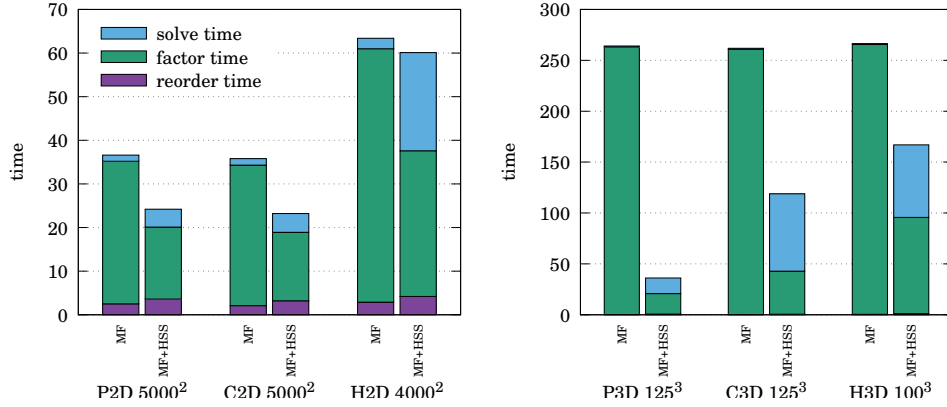


FIG. 8. Summary of the results from Table 4: Poisson (P), convection-diffusion (C), and Helmholtz (H) on 2D (left) and 3D (right) regular grids on a 12-core Intel Ivy Bridge. Poisson and convection-diffusion are in double precision, Helmholtz in complex double precision.

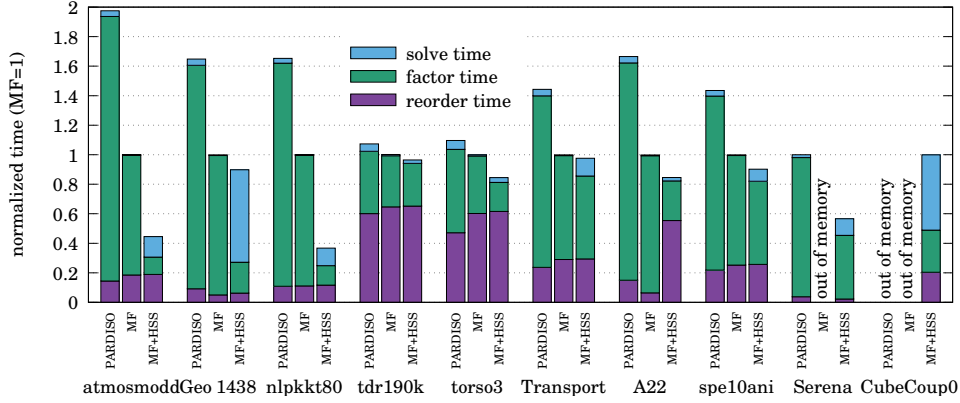


FIG. 9. Comparison of timings for matrices from various applications on a 12-core Intel Ivy Bridge. PARDISO is the sparse direct multithreaded solver from Intel MKL. MF refers to our implementation of the multifrontal method and MF+HSS is our new multifrontal solver with HSS compression. The matrices are taken from the Florida Sparse Matrix Collection and from SciDAC projects at the DOE. For these matrices, which are all quite large and from 2D/3D PDE problems, our MF solver is faster than PARDISO and HSS compression gives an additional speedup.

compared to the standard multifrontal solver. Table 4 also lists the memory usage for each of the experiments, with a maximum reduction in memory usage of a factor 4, i.e., 24% of pure multifrontal.

The times for symbolic factorization in Table 4 are larger for the multifrontal method than for the HSS solver. This is because the memory for dense frontal matrices is allocated during the symbolic factorization while memory for the HSS generators is allocated during the numerical factorization since HSS-ranks are not known in advance.

**6.2. Matrices from various applications.** Figure 9 shows a comparison of timings to solve linear systems with a number of matrices from applications. The matrices ATMOSMODD, GEO\_1438, NLPKKT80, TORSO3, TRANSPORT, and SERENA

TABLE 5

Same as in Figure 9: comparison of timings for matrices from various applications on a 12-core Intel Ivy Bridge. This table also shows memory usage, the optimal number of HSS levels  $\ell_s$ , the optimal compression tolerance  $\varepsilon$ , and the corresponding HSS-rank and number of GMRES(30) iterations. For the SERENA and CUBE\_COUP\_DT0 matrices the pure multifrontal method ran out of memory in double precision.

				MF			HSS						
Matrix	PDE	Order	#nnz	Fact	Solve	Mem	$\ell_s/\ell_{\max}$	$\varepsilon$	Rank	Its	Fact	Solve	Mem
atmosmodd	3D	1.2M	8.8M	81s	0.4s	16GB	6/18	0.9	17	88	25s	11s	5.1GB
Geo_1438	3D	1.4M	63M	205s	1s	40GB	6/18	0.9	8	318	56s	129s	18GB
nlpkkt80	3D	1.1M	29M	197s	0.7s	30GB	6/18	0.5	59	90	49s	23s	12GB
tdr190k	3D	1.1M	43M	19s	0.2s	5.6GB	1/18	$10^{-4}$	61	2	18s	0.4s	5.6GB
torso3	3D	.25M	4.4M	6s	0.05s	1.8GB	6/15	0.5	36	7	5s	0.2s	1.0GB
Transport	3D	1.6M	23M	80s	0.5s	20GB	3/18	$10^{-2}$	182	24	69s	10s	18GB
A22	2D	.59M	145M	127s	0.7s	28GB	10/17	0.1	172	18	105s	3s	6.2GB
spe10-aniso	3D	1.2M	31M	88s	0.4s	19GB	3/18	$10^{-2}$	245	21	73s	7.3s	15GB
Serena*	3D	1.4M	65M	171s	0.5s	22GB	6/18	0.9	11	111	40s	22s	8.1GB
Cube_Coup_dt0*	3D	2.2M	129M	-	-	-	8/19	0.5	100	200	60s	63s	13GB

\*single precision experiment

are from the University of Florida Sparse Matrix Collection.<sup>7</sup> The other matrices, TDR190K, A22, and SPE10-ANISOTROPIC, are from SciDAC projects at the DOE. The matrices are also listed in Table 5, which additionally contains matrix CUBE\_COUP\_DT0. This last matrix is not shown in Figure 9 because our multifrontal code ran out of memory during factorization unless HSS compression was used. All selected matrices are relatively large and originated from a 2D or 3D PDE (on arbitrary domains). In Figure 9, our HSS enabled multifrontal solver (MF+HSS) is compared to the pure multifrontal method (MF) and to the state-of-the-art PARDISO solver [49]. PARDISO, a multithreaded supernodal solver, is part of Intel MKL. For MF and MF+HSS, reorder time includes nested dissection, MC64, and symmetrization of the sparsity structure and for MF+HSS also separator reordering. Factor time includes both symbolic and numerical factorization. The times are normalized to a total time 1 for MF. For the matrices selected for Figure 9, we see a consistent speedup from MF+HSS compared to pure MF, and our MF solver always outperforms the PARDISO solver. PARDISO uses the same METIS nested dissection reordering as our implementation, with comparable reordering times for the different solvers. The supernodal pivoting scheme used in PARDISO for numerical stability does not affect the fill-in, so the overall number of nonzeros in the factors with PARDISO and with our multifrontal code are very similar. Only for the A22 problem does reordering the separator in order to reduce the HSS-ranks take a lot of time. This is probably due to the addition of link-two edges to the separator graph (see section 4.2) since the original matrix already has 246 nonzeros per row on average. However, if those extra edges are not taken into account, the HSS-ranks are much larger, and there is no net performance benefit from using HSS.

**6.3. Many-core parallel performance.** Figure 10 shows performance and parallel scalability of the MF+HSS solver applied to the TORSO3.MTX matrix ( $\ell_s = 6$ ,  $\varepsilon = 0.5$ ) on two leading multicore architectures: a two socket machine with a 12-core Intel Ivy Bridge Xeon per socket and a 60-core Intel Xeon Phi Knight's Corner. When running 12 or fewer threads on the dual socket 24-core Xeon system, the threads are

<sup>7</sup><http://www.cise.ufl.edu/research/sparse/matrices/>

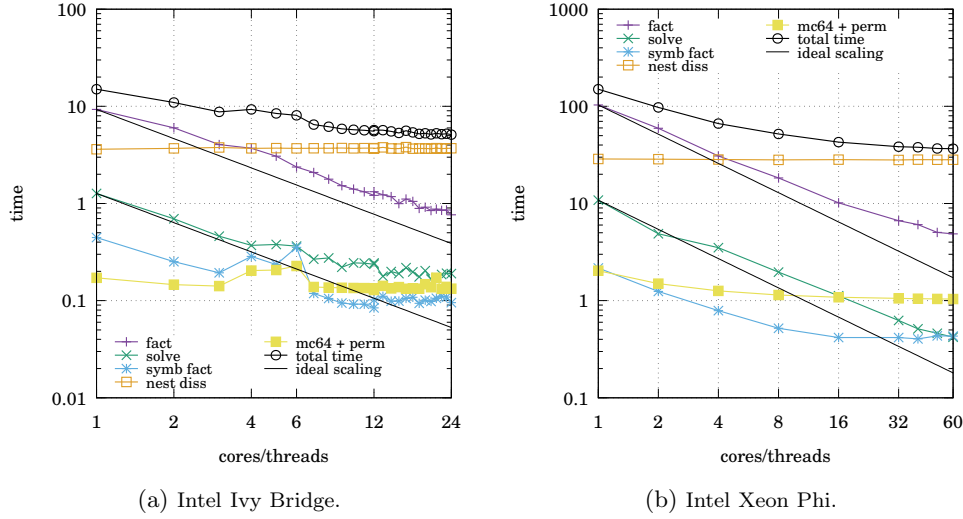


FIG. 10. Multicore scalability of the different steps in the MF+HSS solver on two leading architectures. The MF+HSS solver is applied to the relatively small `TORSO3.MTX` matrix. The code achieves good speedup for the numerical factorization phase and reasonable speedup for the solve (MF+HSS preconditioned GMRES). Note that the sequential reordering codes METIS and MC64 become bottlenecks.

all bound to a single socket (NUMA node). Note that since the Xeon Phi only has 8GB of memory, the larger problems from Table 5 do not fit in its memory. Our code shows good parallel scalability on both architectures for the numerical factorization phase and reasonable scalability for the solve phase. However, with an increasing number of threads the reordering codes MC64 and METIS/SCOTCH quickly become scaling bottlenecks. The MC64 phase in Figure 10 shows some parallel speedup since this time it also includes applying the column permutation from MC64, which is done in parallel.

**7. Conclusions and outlook.** We presented an initial attempt to create a high performance implementation of a novel multifrontal solver with HSS low-rank structures. We show speedups of up to  $7\times$  over the pure multifrontal algorithm for a range of applications. Moreover, our implementation compares favorably to the commercial PARDISO solver. We observed that the new solver has lower computational complexity than the pure multifrontal method. However, the constants involved are much larger. We will focus our attention on trying to reduce these constants (for instance, by trying to reduce the HSS-ranks) and on solving larger problems with a distributed memory implementation. As possible strategies to reduce the HSS-ranks, we consider the following. A power iteration on the random vectors (for instance,  $S^r = (AA^*)^q AR$  with  $q$  a small integer) will improve the quality of the samples at the expense of additional computations; see [29] for further details. We believe the separator reordering (see section 4.2) can be improved, perhaps by taking into account the matrix entries and/or the underlying geometry, also leading to lower ranks. Finally, a better rank-revealing factorization, like a *strong* rank-revealing QR [28], might lead to lower ranks and possibly more stable ULV factorization but at an increased computational cost. The solver with HSS compression achieves lower floating point operation throughput than the pure multifrontal code. Hence, we believe there

is some room for improvement. We will continue performance tuning of the code on various modern computer architectures.

The presented code is part of a package called STRUMPACK. At the moment STRUMPACK has a sparse shared memory solver and a dense distributed memory solver. The longer term goal is to develop and maintain a single scalable code for both sparse and dense problems using hybrid parallelism. The current paper, together with the distributed HSS code developed for [47], represents a good step towards reaching that goal.

The research on fast sparse and dense direct solvers is a very active field at the moment. Some newer algorithmic ideas are, for instance, nested HSS approximation and matrix-free direct solver based preconditioners. In nested HSS approximation, the HSS generators of the frontal matrices are themselves HSS matrices. This could further reduce the overall complexity of the solver. A matrix-free direct solver based preconditioner could be constructed using randomization techniques. It seems that knowledge of the sparsity pattern would be required for this.

**Acknowledgments.** We would like to thank Jianlin Xia for insightful discussions and for his pioneering work on these exciting algorithms. Also thanks to Alex Druinsky for testing an early version of the code. Yvan Notay provided the code to generate the convection diffusion matrices, and the code for the Helmholtz problems was provided by Shen Wang.

#### REFERENCES

- [1] E. AGULLO, A. BUTTARI, A. GUERMOUCHE, AND F. LOPEZ, *Multifrontal QR factorization for multicore architectures over runtime systems*, in Euro-Par 2013 Parallel Processing, Lecture Notes in Comput. Sci. 8097, Springer-Verlag, Berlin, Heidelberg, 2013, pp. 521–532.
- [2] E. AGULLO, J. DEMMEL, J. DONGARRA, B. HADRI, J. KURZAK, J. LANGOU, H. LTAIEF, P. LUSZCZEK, AND S. TOMOV, *Numerical linear algebra on emerging architectures: The plasma and magma projects*, J. Phys. Conf. Ser., 180 (2009), 012037.
- [3] S. AMBIKASARAN, *Fast Algorithms for Dense Numerical Linear Algebra and Applications*, Ph.D. thesis, Stanford, Palo Alto, CA, 2013.
- [4] P. R. AMESTOY, C. ASHCRAFT, O. BOITEAU, A. BUTTARI, J.-Y. L'EXCELLENT, AND C. WEISBECKER, *Improving multifrontal methods by means of block low-rank representations*, SIAM J. Sci. Comput., 37 (2015), pp. A1451–A1474, doi:10.1137/120903476.
- [5] P. R. AMESTOY, T. A. DAVIS, AND I. S. DUFF, *An approximate minimum degree ordering algorithm*, SIAM J. Matrix Anal. Appl., 17 (1996), pp. 886–905, doi:10.1137/S0895479894278952.
- [6] P. R. AMESTOY, I. S. DUFF, J.-Y. L'EXCELLENT, AND J. KOSTER, *A fully asynchronous multifrontal solver using distributed dynamic scheduling*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 15–41, doi:10.1137/S0895479899358194.
- [7] P. R. AMESTOY, I. S. DUFF, D. RUIZ, AND B. UÇAR, *A parallel matrix scaling algorithm*, in High Performance Computing for Computational Science-VECPAR 2008, Lecture Notes in Comput. Sci. 5336, Springer, Berlin, Heidelberg, 2008, pp. 301–313.
- [8] A. AMINFAR, S. AMBIKASARAN, AND E. DARVE, *A fast block low-rank dense solver with applications to finite-element matrices*, J. Comput. Phys., 304 (2016), pp. 170–188.
- [9] C. AUGONNET, S. THIBAUT, R. NAMYST, AND P.-A. WACRENIER, *StarPU: A unified platform for task scheduling on heterogeneous multicore architectures*, Concurr. Comput., 23 (2011), pp. 187–198.
- [10] M. BEBENDORF, *Approximation of boundary element matrices*, Numer. Math., 86 (2000), pp. 565–589.
- [11] R. D. BLUMOFFE, C. F. JOERG, B. C. KUSZMAUL, C. E. LEISERSON, K. H. RANDALL, AND Y. ZHOU, *Cilk: An efficient multithreaded runtime system*, in PPOPP '95, Proceedings of the Fifth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, ACM, New York, 1996, pp. 207–216.

- [12] R. D. BLUMOFE AND C. E. LEISERSON, *Scheduling multithreaded computations by work stealing*, J. ACM, 46 (1999), pp. 720–748.
- [13] S. BÖRM, L. GRASEDYCK, AND W. HACKBUSCH, *Introduction to hierarchical matrices with applications*, Eng. Anal. Bound. Elem., 27 (2003), pp. 405–422.
- [14] G. BOSILCA, A. BOUTELLER, A. DANALIS, T. HERAULT, P. LEMARINIER, AND J. DONGARRA, *Dague: A generic distributed dag engine for high performance computing*, Parallel Comput., 38 (2012), pp. 37–51.
- [15] A. BUTTARI, J. LANGOU, J. KURZAK, AND J. DONGARRA, *A class of parallel tiled linear algebra algorithms for multicore architectures*, Parallel Comput., 35 (2009), pp. 38–53.
- [16] T. F. CHAN, *Rank revealing QR factorizations*, Linear Algebra Appl., 88 (1987), pp. 67–82.
- [17] S. CHANDRASEKARAN, P. DEWILDE, M. GU, AND N. SOMASUNDERAM, *On the numerical rank of the off-diagonal blocks of Schur complements of discretized elliptic PDEs*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2261–2290, doi:10.1137/090775932.
- [18] S. CHANDRASEKARAN, M. GU, AND T. PALS, *A fast ULV decomposition solver for hierarchically semiseparable representations*, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 603–622, doi:10.1137/S0895479803436652.
- [19] H. CHENG, Z. GIMBUTAS, P.-G. MARTINSSON, AND V. ROKHLIN, *On the compression of low rank matrices*, SIAM J. Sci. Comput., 26 (2005), pp. 1389–1404, doi:10.1137/030602678.
- [20] A. R. CURTIS AND J. K. REID, *On the automatic scaling of matrices for Gaussian elimination*, J. Inst. Math. Appl., 10 (1972), pp. 118–124.
- [21] J. DONGARRA, M. FAVERGE, H. LTAIEF, AND P. LUSZCZEK, *Achieving numerical accuracy and high performance using recursive tile lu factorization with partial pivoting*, Concurr. Comput., 26 (2014), pp. 1408–1431.
- [22] I. S. DUFF AND J. KOSTER, *The design and use of algorithms for permuting large entries to the diagonal of sparse matrices*, SIAM J. Matrix Anal. Appl., 20 (1999), pp. 889–901, doi:10.1137/S0895479897317661.
- [23] I. S. DUFF AND J. K. REID, *The multifrontal solution of indefinite sparse symmetric linear*, ACM Trans. Math. Software, 9 (1983), pp. 302–325.
- [24] A. DURAN, E. AYGUADÉ, R. M. BADIA, J. LABARTA, L. MARTINELL, X. MARTORELL, AND J. PLANAS, *Ompss: A proposal for programming heterogeneous multi-core architectures*, Parallel Process. Lett., 21 (2011), pp. 173–193.
- [25] B. ENGQUIST AND L. YING, *Sweeping preconditioner for the Helmholtz equation: Hierarchical matrix representation*, Comm. Pure Appl. Math., 64 (2011), pp. 697–735.
- [26] M. FRIGO, C. E. LEISERSON, H. PROKOP, AND S. RAMACHANDRAN, *Cache-oblivious algorithms*, in 40th Annual Symposium on Foundations of Computer Science, IEEE, Washington, DC, 1999, pp. 285–297.
- [27] S. GHEMAWAT AND P. MENAGE, *Tcmalloc: Thread-Caching Malloc*, <http://goog-perftools.sourceforge.net/doc/tcmalloc.html> (2009).
- [28] M. GU AND S. C. EISENSTAT, *Efficient algorithms for computing a strong rank-revealing QR factorization*, SIAM J. Sci. Comput., 17 (1996), pp. 848–869, doi:10.1137/0917055.
- [29] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53 (2011), pp. 217–288, doi:10.1137/090771806.
- [30] P. HÉNON, P. RAMET, AND J. ROMAN, *PaStiX: A high-performance parallel direct solver for sparse symmetric positive definite systems*, Parallel Comput., 28 (2002), pp. 301–321.
- [31] J. D. HOGG, J. K. REID, AND J. A. SCOTT, *Design of a multicore sparse Cholesky factorization using DAGs*, SIAM J. Sci. Comput., 32 (2010), pp. 3627–3649, doi:10.1137/090757216.
- [32] G. KARYPIS AND V. KUMAR, *A fast and high quality multilevel scheme for partitioning irregular graphs*, SIAM J. Sci. Comput., 20 (1998), pp. 359–392, doi:10.1137/S1064827595287997.
- [33] K. KIM AND V. EIJKHOUT, *A parallel sparse direct solver via hierarchical DAG scheduling*, ACM Trans. Math. Software, 41 (2014), 3.
- [34] R. KRIEMANN, *H-LU factorization on many-core systems*, Comput. Vis. Sci., 16 (2013), pp. 105–117.
- [35] X. LACOSTE, M. FAVERGE, G. BOSILCA, P. RAMET, AND S. THIBAUT, *Taking advantage of hybrid systems for sparse direct solvers via task-based runtimes*, in 2014 IEEE International Parallel & Distributed Processing Symposium Workshops (IPDPSW), IEEE, Washington, DC, 2014, pp. 29–38.
- [36] J.-Y. L’EXCELLENT, *Multifrontal Methods: Parallelism, Memory Usage and Numerical Aspects*, Habilitation à Diriger des Recherches, École Normale Supérieure de Lyon, Lyon, France, 2012.
- [37] L. LIN, J. LU, AND L. YING, *Fast construction of hierarchical matrix representation from matrix-vector multiplication*, J. Comput. Phys., 230 (2011), pp. 4071–4087.

- [38] J. W. H. LIU, *The multifrontal method for sparse matrix solution: Theory and practice*, SIAM Rev., 34 (1992), pp. 82–109, doi:10.1137/1034004.
- [39] P.-G. MARTINSSON, *A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1251–1274, doi:10.1137/100786617.
- [40] J. D. MCCALPIN, *STREAM: Sustainable Memory Bandwidth in High Performance Computers*, <http://www.cs.virginia.edu/stream/>.
- [41] M. MCCOOL, J. REINDERS, AND A. ROBISON, *Structured Parallel Programming: Patterns for Efficient Computation*, Morgan Kaufmann, San Francisco, CA, 2012.
- [42] A. NAPOV AND X. S. LI, *An algebraic multifrontal preconditioner that exploits the low-rank property*, Numer. Linear Algebra Appl., 23 (2016), pp. 61–82, doi:10.1002/nla.2006.
- [43] Y. NOTAY, *An aggregation-based algebraic multigrid method*, Electron. Trans. Numer. Anal., 37 (2010), pp. 123–146.
- [44] F. PELLEGRINI AND J. ROMAN, *Scotch: A software package for static mapping by dual recursive bipartitioning of process and architecture graphs*, in High-Performance Computing and Networking, Springer-Verlag, London, 1996, pp. 493–498.
- [45] G. QUINTANA-ORTÍ, X. SUN, AND C. H. BISCHOF, *A BLAS-3 version of the QR factorization with column pivoting*, SIAM J. Sci. Comput., 19 (1998), pp. 1486–1494, doi:10.1137/S1064827595296732.
- [46] J. REINDERS, *Intel Threading Building Blocks: Outfitting C++ for Multi-core Processor Parallelism*, O'Reilly Media, Sebastopol, CA, 2007.
- [47] F.-H. ROUET, X. S. LI, P. GHYSELS, AND A. NAPOV, *A distributed-memory package for dense hierarchically semi-separable matrix computations using randomization*, ACM Trans. Math. Software, to appear.
- [48] D. RUIZ, *A Scaling Algorithm to Equilibrate Both Rows and Columns Norms in Matrices*, Technical report RT/APO/01/4, ENSEEIHT-IRIT, Toulouse, France, 2001; also appeared as RAL report RAL-TR-2001-034, 2001.
- [49] O. SCHENK, K. GÄRTNER, AND W. FICHTNER, *Efficient sparse lu factorization with left-right looking strategy on shared memory multiprocessors*, BIT, 40 (2000), pp. 158–176.
- [50] R. VANDEBRIL, M. VAN BAREL, G. GOLUB, AND N. MASTRONARDI, *A bibliography on semiseparable matrices*, Calcolo, 42 (2005), pp. 249–270.
- [51] S. WANG, M. V. DE HOOP, AND J. XIA, *On 3D modeling of seismic wave propagation via a structured parallel multifrontal direct Helmholtz solver*, Geophys. Prospect., 59 (2011), pp. 857–873.
- [52] S. WANG, X. S. LI, F.-H. ROUET, J. XIA, AND M. V. DE HOOP, *A parallel geometric multifrontal solver using hierarchically semiseparable structure*, ACM Trans. Math. Software, 42 (2016), 21.
- [53] C. WEISBECKER, *Improving Multifrontal Solvers by Means of Algebraic Block Low-Rank Representations*, Ph.D. thesis, Institut National Polytechnique de Toulouse, Toulouse, France, 2013.
- [54] J. H. WILKINSON, *Rounding Errors in Algebraic Processes*, Dover Publications, New York, 1994.
- [55] J. XIA, *Randomized sparse direct solvers*, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 197–227, doi:10.1137/12087116X.
- [56] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Fast algorithms for hierarchically semiseparable matrices*, Numer. Linear Algebra Appl., 17 (2010), pp. 953–976.
- [57] J. XIA, Y. XI, AND M. GU, *A superfast structured solver for Toeplitz linear systems via randomized sampling*, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 837–858, doi:10.1137/12087116X.
- [58] A. YARKHAN, J. KURZAK, AND J. DONGARRA, *Quark Users Guide: Queueing and Runtime for Kernels*, University of Tennessee Innovative Computing Laboratory Technical Report ICL-UT-11-02, University of Tennessee, Knoxville, TN, 2011.